

# Competitive Minimax Universal Decoding for Several Ensembles of Random Codes \*

Yaniv Akirav and Neri Merhav  
 Department of Electrical Engineering  
 Technion - Israeli Institute of Technology  
 Technion City, Haifa 32000, Israel  
 Emails:[yaniva@tx,merhav@ee].technion.ac.il

February 1, 2008

## Abstract

Universally achievable error exponents pertaining to certain families of channels (most notably, discrete memoryless channels (DMC's)), and various ensembles of random codes, are studied by combining the competitive minimax approach, proposed by Feder and Merhav, with Chernoff bound and Gallager's techniques for the analysis of error exponents. In particular, we derive a single-letter expression for the largest, universally achievable fraction  $\xi$  of the optimum error exponent pertaining to the optimum ML decoding. Moreover, a simpler single-letter expression for a lower bound to  $\xi$  is presented. To demonstrate the tightness of this lower bound, we use it to show that  $\xi = 1$ , for the binary symmetric channel (BSC), when the random coding distribution is uniform over: (i) all codes (of a given rate), and (ii) all linear codes, in agreement with well-known results. We also show that  $\xi = 1$  for the uniform ensemble of systematic linear codes, and for that of time-varying convolutional codes in the bit-error-rate sense. For the latter case, we also show how the corresponding universal decoder can be efficiently implemented using a slightly modified version of the Viterbi algorithm which employs two trellises.

**Index Terms:** error exponent, universal decoding, generalized likelihood ratio test, channel uncertainty, competitive minimax, Viterbi algorithm, maximum mutual information decoding.

---

\*This research was supported by the Israel Science Foundation (ISF), grant no. 223/05.

# 1 Introduction

In many real-life situations, encountered in digital coded communication systems, channel variability and uncertainty prohibit the use of the optimum maximum likelihood (ML) decoder, and so, universal decoders, independent of the unknown channel parameters, are sought.

The topic of universal coding and decoding for unknown channels has received considerable attention in the last three decades. In [5], Goppa offered the *maximum mutual information* (MMI) decoder, which decides in favor of the code vector with maximum empirical mutual information with the channel output. Goppa showed that for DMC's, MMI decoding achieves capacity. Csiszár and Körner [2] also explored the universal decoding problem for DMC's with finite input and output alphabet. They showed that the random coding error exponent associated with a uniform random coding distribution over a type class achieves the optimum error exponent. Csiszár [1] proved that for any channel within the class of DMC's with additive noise, and the uniform random coding distribution over linear codes, the optimum error exponent is achievable by a decoder minimizing the noise empirical entropy, universally for all the channels in the class. Ziv [12] explored the universal decoding problem for finite state channels with finite input and output alphabets, for which the next channel state is a deterministic (but unknown) function of the channel current state and current inputs and outputs. For codes governed by a uniform random coding over a given set, he proved that a decoder based on the Lempel-Ziv algorithm asymptotically achieves the error exponent associated with ML decoding. In [6], Ziv and Lapidot proved that the latter decoder is universal for a wider class of finite-state channels. In [3], Feder and Lapidot found sufficient conditions for families of channels, to have universal decoders that asymptotically achieve the random coding error exponent associated with ML decoding.

Universal coding and decoding were explored also with regard to the generalized likelihood ratio test (GLRT). In this approach, each message is scored according to the maximum likelihood (over the parameter space) of the channel output vector given the message, and a decision is made in favor of the message that attains the highest maximum likelihood. Although provably optimum in certain asymptotic situations [11], [2, p. 165, Theorem 5.2], there are cases where the GLRT is strictly suboptimum [6, Sect. III, pp. 1754–1755], [4,

Appendix].

The competitive minimax criterion, first presented in [4], is an attempt for a general methodological approach to the problem of universal decoding. According to this approach, the criterion is the minimum (over all decision rules) of the maximum (over all channels in the family) of the ratio between the error probability associated with a given channel and given decision rule, and the error probability of the ML decoder for that channel, raised to some power  $\xi \in [0, 1]$  (cf. eq. (2) below). The largest power  $\xi = \xi^*$  such that the value of this minimax ratio does not grow exponentially with the block length, is the maximum universally achievable fraction of the ML error exponent.

The main contribution of this paper is in deriving a single-letter expression to  $\xi^*$ , in terms of the rate  $R$  and a general random coding distribution, for fairly general families of channels and ensembles of random codes. While in previous works the universality was proved for certain channel models (e.g. finite-state channels, etc.) and random coding distributions (e.g. uniform distribution over a given type class, etc.), this work deals with general families of DMC's (cf. Sect. II) and general random coding distributions (cf. eq. (7)). We should note that a similar technique can be used to broaden the result for  $\xi^*$  to other channel families, e.g. Markov channels, finite state channels, etc.

In addition, a single-letter expression for a lower bound to  $\xi^*$  is presented, which is simpler to work with, and is believed to be tight. This lower bound is true also for random coding distribution over ensembles of linear code and systematic linear codes. The tightness of this lower bound is demonstrated for the case of the BSC. For this model, we show that  $\xi^* = 1$ , when the random coding distribution is uniform over all codes and over all linear codes, in agreement with well-known results. We also show that  $\xi^* = 1$  for the ensemble of systematic linear codes, and for that of time-varying convolutional codes in the bit-error-rate sense. Using the fact that in the case of the BSC, the minimax decoding metric degenerates to a simpler metric, we propose an efficient implementation based on a slightly modified version of the Viterbi algorithm.

The outline of the paper is as follows. In Section II, we establish the notation that will be used throughout the paper and provide a formal definition of the universal decoding problem. In Section III, the main results are stated and discussed. Section IV contains a detailed proof of the single-letter expression for  $\xi^*$  will be provided. In Section V, the tightness of the lower bound to  $\xi^*$  is demonstrated for the case of the BSC with an unknown

crossover probability. In Section VI, we prove that for the ensemble of time-varying convolutional codes and the BSC with an unknown crossover probability, the minimax decoder achieves the same bit error exponent as the ML decoder, which is used when the parameter is known.

## 2 Notation and Problem Definition

Throughout this paper, scalar random variables (RV's) will be denoted by capital letters, their sample values will be denoted by the respective lower case letters, and their alphabets will be denoted by the respective calligraphic letters. A similar convention will apply to random vectors of dimension  $N$  and their sample values, which will be denoted with same symbols in the bold face font. The set of all  $N$ -vectors with components taking values in a certain alphabet, will be denoted as the same alphabet superscripted by  $N$ .

Information theoretic quantities like entropies, conditional entropies, and mutual informations, will be denoted following the usual conventions of the information theory literature, e.g.,  $H(X)$ ,  $H(X|Y)$ ,  $I(X;Y)$ , and so on. With a slight abuse of notation, when we wish to emphasize the dependence of the entropy on the underlying probability distribution  $P$ , we denote it by  $H(P)$ .

The mutual information between the input and the output of the channel

$\{P_\theta(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$ , when the input is governed by  $Q$ , will be denoted by

$$I_\theta(Q) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} Q(x) P_\theta(y|x) \ln \frac{P_\theta(y|x)}{\sum_{x \in \mathcal{X}} Q(x) P_\theta(y|x)}, \quad (1)$$

and the capacity of the channel will be denoted by  $C_\theta = \max_Q I_\theta(Q)$ .

The number of occurrences of a letter  $a \in \mathcal{X}$  in a vector  $\mathbf{x} \in \mathcal{X}^N$  will be denoted by  $N_{\mathbf{x}}(a)$ . The empirical distribution of  $\mathbf{x}$  will be denoted by  $P_{\mathbf{x}} = \{P_{\mathbf{x}}(a) = N_{\mathbf{x}}(a)/N, a \in \mathcal{X}\}$ . The type class of  $\mathbf{x}$  is defined as  $T_{\mathbf{x}} = \{\mathbf{x}' : P_{\mathbf{x}'} = P_{\mathbf{x}}\}$  and  $H_{\mathbf{x}}(X) = -\sum_{a \in \mathcal{X}} P_{\mathbf{x}}(a) \ln P_{\mathbf{x}}(a)$  will denote the entropy of a random variable (RV)  $X$ , with distribution  $P_{\mathbf{x}}$ . Similarly, the number of occurrences of a letter pair  $(a, b) \in \mathcal{X} \times \mathcal{Y}$  in the vector pair  $(\mathbf{x}, \mathbf{y})$  will be denoted by  $N_{\mathbf{xy}}(a, b)$ ,  $P_{\mathbf{xy}} = \{P_{\mathbf{xy}}(a, b) = N_{\mathbf{xy}}(a, b)/N, (a, b) \in \mathcal{X} \times \mathcal{Y}\}$  will denote the joint empirical distribution of  $(\mathbf{x}, \mathbf{y})$ ,  $T_{\mathbf{xy}} = \{\mathbf{x}', \mathbf{y}' : P_{\mathbf{x}'\mathbf{y}'} = P_{\mathbf{xy}}\}$  will stand for the joint type class of  $(\mathbf{x}, \mathbf{y})$ , and  $H_{\mathbf{xy}}(X, Y) = -\sum_{a, b \in \mathcal{X} \times \mathcal{Y}} P_{\mathbf{xy}}(a, b) \ln P_{\mathbf{xy}}(a, b)$  will denote the joint entropy of RV's  $(X, Y)$  with joint distribution  $P_{\mathbf{xy}}$ . We will use  $T_{\mathbf{x}|\mathbf{y}} = \{\mathbf{x}' : P_{\mathbf{x}'\mathbf{y}} = P_{\mathbf{xy}}\}$  to denote the conditional type class of  $\mathbf{x}$  given  $\mathbf{y}$ ,  $P_{\mathbf{x}|\mathbf{y}}(a|b) = N_{\mathbf{xy}}(a, b)/N_{\mathbf{y}}(b)$ ,  $(a, b) \in$

$\mathcal{X} \times \mathcal{Y}$ , to denote the conditional empirical distribution related to  $(a, b) \in \mathcal{X} \times \mathcal{Y}$ , and  $H_{\mathbf{x}\mathbf{y}}(X|Y) = -\sum_{a,b \in \mathcal{X} \times \mathcal{Y}} P_{\mathbf{x}\mathbf{y}}(a, b) \ln P_{\mathbf{x}|\mathbf{y}}(a|b)$  to denote the conditional entropy of  $X$  given  $Y$ , induced by the joint distribution  $P_{\mathbf{x}\mathbf{y}}$ . The empirical mutual information between RV's  $X$  and  $Y$  with joint distribution  $P_{\mathbf{x}\mathbf{y}}$  will be denoted by  $I_{\mathbf{x}\mathbf{y}}(X; Y) = H_{\mathbf{x}}(X) - H_{\mathbf{x}\mathbf{y}}(X|Y)$ .

The expectation of a function  $F(X, Y)$ , where  $X$  and  $Y$  are RV's distributed according to the empirical distribution of  $\mathbf{x}$  and  $\mathbf{y}$ , will be denoted by

$$\hat{E}_{\mathbf{x}\mathbf{y}}\{F(X, Y)\} = \sum_{a \in \mathcal{X}} \sum_{b \in \mathcal{Y}} P_{\mathbf{x}\mathbf{y}}(a, b) F(a, b).$$

The notation  $E_Q\{F(\mathbf{X})\}$  will be used for the expectation of a function  $F(\mathbf{X})$ , where the random vector  $\mathbf{X}$  is governed by  $Q$ .

The Hamming distance between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  will be denoted by  $d(\mathbf{x}, \mathbf{y})$ , and its normalization by  $N$  will be denoted by  $\delta(\mathbf{x}, \mathbf{y})$ . For a finite set  $\mathcal{A}$ ,  $|\mathcal{A}|$  will stand for its cardinality. The divergence between two probability measures  $P$  and  $Q$  over an alphabet  $\mathcal{U}$  will be denoted by  $D(P||Q) = \sum_{u \in \mathcal{U}} P(u) \ln \frac{P(u)}{Q(u)}$ , where  $0 \ln 0$  and  $0 \ln \frac{0}{0}$  are defined as 0, and  $P \ln \frac{P}{0}$  for  $P > 0$  is defined as  $\infty$ . For two positive sequences  $\{A_N\}_{N \geq 1}$  and  $\{B_N\}_{N \geq 1}$ , the notation  $A_N \doteq B_N$  will express the fact that  $\{A_N\}_{N \geq 1}$  and  $\{B_N\}_{N \geq 1}$  are of the same exponential order, i.e.,

$$\lim_{N \rightarrow \infty} \frac{1}{N} \ln (A_N/B_N) = 0.$$

Consider a DMC with a finite input alphabet  $\mathcal{X}$ , a finite output alphabet  $\mathcal{Y}$ , and single letter transition probabilities  $\{P_\theta(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$ , where  $\theta$  is an unknown parameter vector, taking values in some set  $\Theta$ . The channel is fed by an input vector of length  $N$ ,  $\mathbf{x} \in \mathcal{X}^N$ , and generates an output vector  $\mathbf{y} \in \mathcal{Y}^N$  according to  $P_\theta(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^N P_\theta(y_i|x_i)$ . A rate- $R$  block code of length  $N$  consists of  $M = e^{NR}$   $N$ -vectors  $\mathbf{x}_m \in \mathcal{X}^N$ ,  $0 \leq m \leq M-1$ , representing  $M$  different messages. A decoder  $\Omega$  is a partition of  $\mathcal{Y}^N$  into  $M$  regions,  $\Omega_0, \Omega_1, \dots, \Omega_{M-1}$ , such that if  $\mathbf{y}$  falls into  $\Omega_m$ , a decision is made in favor of message  $m$ .

Given a code  $\mathcal{C}$ , the competitive minimax criterion [4] is defined as

$$S_N \triangleq \min_{\Omega} \max_{\theta \in \Theta} \left\{ \frac{P_E(\Omega|\theta)}{[P_E^*(\theta)]^\xi} \right\}, \quad 0 \leq \xi \leq 1, \quad (2)$$

where  $P_E(\Omega|\theta) = \frac{1}{M} \sum_{m=0}^{M-1} \sum_{\mathbf{y} \in \Omega_m^c} P_\theta(\mathbf{y}|\mathbf{x}_m)$  is the error probability related to a decoder  $\Omega$  for a given value of  $\theta$ , and  $P_E^*(\theta) = \min_{\Omega} P_E(\Omega|\theta)$  is the ML decoding error probability when  $\theta$  is known.

The ratio  $P_E(\Omega|\theta)/[P_E^*(\theta)]^\xi$  designates the loss in error probability, caused by using a universal decoder which is ignorant of  $\theta$ , relative to the optimal ML decoding for that  $\theta$ . The parameter  $\xi$  can be interpreted as the fraction of the optimal error exponent to which the universal decoder error exponent is compared. In order to minimize this loss uniformly over all  $\Theta$ , a decoder  $\Omega$  which minimizes the worst case of that ratio (i.e., its maximum), is sought.

As  $S_N$  addresses the ratio between the error probabilities, it corresponds to the difference between the error exponents related to these errors. It is well known that for most channels, the decoding error decays exponentially with the block length  $N$ . Therefore, if the value of  $S_N$ , for a decision rule  $\Omega$  achieved by (2), grows sub-exponentially with  $N$ , i.e.,  $\lim_{N \rightarrow \infty} \frac{1}{N} \ln S_N = 0$ , it means that, uniformly over  $\Theta$ , the error probability associated with  $\Omega$  decays with an exponential rate which is at least a fraction  $\xi$  of the error exponent rate of  $P_E^*(\theta)$ .

In [4], the following decision rule has been shown to be asymptotically optimal in the minimax sense for a given  $\xi$ :

$$\Omega_m = \{\mathbf{y} | f(\mathbf{x}_m, \mathbf{y}) \geq f(\mathbf{x}_{m'}, \mathbf{y}), \forall m' \neq m\} \quad (3)$$

with ties broken arbitrarily, where

$$f(\mathbf{x}, \mathbf{y}) \triangleq \max_{\theta \in \Theta} f_\theta(\mathbf{x}, \mathbf{y}), \quad (4)$$

$$f_\theta(\mathbf{x}, \mathbf{y}) \triangleq \frac{1}{N} \ln P_\theta(\mathbf{y}|\mathbf{x}) + \xi E^*(\theta), \quad (5)$$

and  $E^*(\theta)$  stands for the asymptotic exponent associated with  $P_E^*(\theta)$ . A decoder  $\Omega$ , defined by (3), will be called the *minimax decoder* hereafter.

A natural question that may arise, at this point, is with regard to the choice of the free parameter  $\xi$ . As mentioned above, the main guideline proposed in [4] is to seek the maximum value  $\xi^*$  of  $\xi$  such that  $S_N$  would still grow sub-exponentially with  $N$ .

In the random coding regime, the error probabilities at the numerator and the denominator of (2) are replaced by the corresponding average error probabilities, i.e.,

$$\overline{S}_N \triangleq \min_{\Omega} \max_{\theta \in \Theta} \left\{ \frac{\overline{P}_E(\Omega|\theta)}{[\overline{P}_E^*(\theta)]^\xi} \right\} \quad (6)$$

and the decoder (3) is used, with  $E^*(\theta)$  being replaced by  $E_r^*(\theta)$ , the random coding error exponent associated with  $\overline{P}_E^*(\theta)$ .

The main purpose of this paper is to translate the above-mentioned guideline for the choice of  $\xi$  into a concrete single-letter formula for the random coding regime.

### 3 Statement of Results

In this section, by evaluating the exponential order of  $\overline{S}_N$ , we derive a formula for  $\xi^*$ , the largest value of  $\xi$  for which  $\overline{S}_N$  is sub-exponential in  $N$ . Moreover, an expression for the lower bound to  $\xi^*$  is also derived, and its tightness is demonstrated for the BSC model and for several ensembles of random codes.

#### 3.1 General codes

We begin with a few definitions. For every positive integer  $N$ , let  $Q_N$  be a random coding distribution for  $N$ -vectors, of the following form:

$$Q_N(\mathbf{x}) = \frac{Q_N(T\mathbf{x})}{|T\mathbf{x}|}, \quad (7)$$

i.e., uniform distribution for all the vectors within the same type class. Of course,

$$\sum_{T\mathbf{x}} Q_N(T\mathbf{x}) = 1.$$

Now, let

$$\Delta_N(P\mathbf{x}) = -\frac{1}{N} \ln Q_N(T\mathbf{x}),$$

and let  $\Delta_N^*(P)$  be an extension of the function  $\Delta_N(P\mathbf{x})$  that is defined over the continuum of probability distributions over  $\mathcal{X}$  (rather than just the set of rational probability distributions with denominator  $N$ ). We next define the class  $\mathcal{Q}$  of sequences of random coding distributions  $\{Q_N\}$  as follows: A sequence of random coding distributions  $\{Q_N\}_{N \geq 1}$  is said to belong to the class  $\mathcal{Q}$  if there exists such an extension  $\Delta_N^*(P)$  that converges, as  $N \rightarrow \infty$ , to a certain non-negative functional  $\Delta^*(P)$ , uniformly over all probability distributions  $\{P\}$  over  $\mathcal{X}$ .

It is easy to see that the class  $\mathcal{Q}$  essentially covers all random coding distributions that are customarily used (and much more). In particular, to approximate a random coding distribution which is uniform within a small neighborhood of one type class – corresponding to a probability distribution  $P_0$ , and which vanishes elsewhere, we set  $\Delta^*(P) = 0$  for every  $P$  in that neighborhood of  $P_0$ , and  $\Delta^*(P) = \infty$  elsewhere. For the case where

$Q$  is i.i.d.,  $\Delta^*(P) = D(P\|Q)$ . In particular, if  $Q(\mathbf{x}) = 1/|\mathcal{X}|^N$  for all  $\mathbf{x} \in \mathcal{X}^N$ , then  $\Delta^*(P) = \ln |\mathcal{X}| - H(P)$ .

Given a joint distribution  $P_{XY}$ , a real  $\alpha$ , and a value of  $\theta \in \Theta$ , let

$$A(\theta, \alpha, P_{XY}) \triangleq I(X; Y) + \Delta^* \left( \sum_{b \in \mathcal{Y}} P_Y(b) P_{X|Y}(\cdot|b) \right) - \alpha E \ln P_\theta(Y|X), \quad (8)$$

where  $E\{\cdot\}$  is the expectation and  $I(X; Y)$  is the mutual information w.r.t. a generic joint distribution  $P_{XY}(a, b) = P_Y(b) P_{X|Y}(a|b)$  of the RV's  $(X, Y)$ .

Next, for distributions  $P_Y, P_{X|Y}$  and  $P_{X'|Y}$ , two parameters  $\theta, \theta' \in \Theta$ , and reals  $0 \leq \rho \leq 1$  and  $s \geq 0$ , define:

$$B(\theta, \theta', P_Y, P_{X|Y}, P_{X'|Y}, s, \rho) \triangleq A(\theta, 1 - s\rho, P_{XY}) + \rho \cdot A(\theta', s, P_{X'Y}) - H(Y), \quad (9)$$

where  $H(Y)$  is the entropy of  $Y$  induced by  $P_Y$ . Finally, let

$$\begin{aligned} \xi^*(R) = \min_{P_{XY}} \min_{\theta' \in \Theta} \max \left\{ \min_{\theta \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \min_{P_{X'|Y}} \frac{B(\theta, \theta', P_Y, P_{X|Y}, P_{X'|Y}, s, \rho) - \rho R}{(1 - \rho s) E_r^*(\theta) + \rho s E_r^*(\theta')}, \right. \\ \left. \max_{\theta \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \min_{P_{X'|Y}} \frac{B(\theta, \theta', P_Y, P_{X|Y}, P_{X'|Y}, s, \rho) - \rho R}{(1 - \rho s) E_r^*(\theta) + \rho s E_r^*(\theta')} \right\} \end{aligned} \quad (10)$$

Our main result, in this section, is the following:

**Theorem 1** *Consider a sequence of ensembles of codes, where each codeword is drawn independently, under a distribution  $Q_N$ , and the sequence  $\{Q_N\}_{N \geq 1}$  is a member of the class  $\mathcal{Q}$ . Then,*

1. *For every  $\xi \leq \xi^*(R)$ ,  $\lim_{N \rightarrow \infty} \frac{1}{N} \ln \bar{S}_N \leq 0$ .*
2. *There exists a sequence of encoders  $\{C_N\}_{N \geq 1}$  and minimax decoders  $\{\Omega_N\}_{N \geq 1}$  with  $\xi = \xi^*(R)$ , for which:*

$$\liminf_{N \rightarrow \infty} \left[ -\frac{1}{N} \ln P_E(\Omega_N|\theta) \right] \geq \xi \cdot E^*(\theta)$$

*uniformly over  $\theta \in \Theta$ .*

3. *For every  $\xi > \xi^*(R)$ ,  $\lim_{N \rightarrow \infty} \frac{1}{N} \ln \bar{S}_N > 0$ .*



The proof of Theorem 1 appears in Section IV.

We now pause to discuss Theorem 1 and some of its aspects.

The theorem suggests a conceptually simple strategy for universal decoding: Given  $R$  and the sequence  $\{Q_N\}_{N \geq 1}$ , first, compute  $\xi^*(R)$  using eq. (10). This may require some non-trivial optimization procedures, but it has to be done only once. It should be mentioned that if closed-form analytic expression does not seem available, the computation can be carried out at least numerically, since this is a single-letter expression. Once  $\xi^*(R)$  has been computed, apply the minimax decoding rule with  $\xi = \xi^*(R)$  and the theorem guarantees that the resulting random coding error exponent associated with the decoder is as specified in the second item of that theorem. Moreover, the third item of the theorem implies that in the random coding regime,  $\xi^*(R)$  is the largest fraction of  $E^*(\theta)$  that is uniformly achievable by a universal decoder.

As mentioned earlier, when  $Q$  is uniform i.i.d.,  $\Delta^*(P) = \ln |\mathcal{X}| - H(X)$  (where  $X$  is governed by  $P$ ), and therefore

$$A(\theta, \alpha, P_{XY}) = \ln |\mathcal{X}| - H(X|Y) - \alpha E \ln P_\theta(Y|X). \quad (11)$$

This observation will be used in Section V which deals with the BSC model, as well as in Section A.1 of the Appendix (ensembles of linear and systematic linear codes), as they both assume a binary i.i.d. random coding distribution.

The theorem is interesting, of course, only when  $\xi^*(R) > 0$ , which is the case in many situations, at least as long as  $R$  is not too large. It should be pointed out that the exponential rate  $\xi^*(R) \cdot E^*(\theta)$ , guaranteed by Theorem 1, is only a lower bound to the real exponential rate (as the minimax criterion is aimed to consider all  $\theta \in \Theta$ ), and that true exponential rate, at some points in  $\Theta$ , might be larger.

As mentioned above, the exact formula for  $\xi^*$ , given in eq. (10), includes many optimizations and hence might be complicated for calculation. Therefore, we next present a simpler expression for a lower bound to  $\xi^*$ , denoted by  $\xi_{LB}^*(R)$ , which we believe is tight at least for several families of channels. Another motivation for presenting  $\xi_{LB}^*(R)$  is that it holds also for ensembles of linear and systematic linear codes, as we will shall in the next subsection. The expression for  $\xi_{LB}^*(R)$  will be derived from  $\xi^*(R)$  by: (i) avoiding the inner maximization between two terms in (10) by choosing the left term, and (ii) interchanging

between the minimization over  $P_{X|Y}$  and the maximization over  $\lambda$  and  $\rho$ , i.e:

$$\xi_{LB}^*(R) \triangleq \min_{P_Y} \min_{\theta \in \Theta} \min_{\theta' \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \min_{P_{X|Y}} \min_{P_{X'|Y}} \frac{B(\theta, \theta', P_Y, P_{X|Y}, P_{X'|Y}, \lambda, \rho) - \rho R}{(1 - \lambda\rho) \cdot E_r^*(\theta) + \lambda\rho \cdot E_r^*(\theta')}. \quad (12)$$

As  $\xi_{LB}^*(R)$  is a lower bound to  $\xi^*$ , it is obvious to see that parts 1 and 2 of Theorem 1 hold for it as well.

### 3.2 Linear codes

We next provide a variation of  $\xi_{LB}^*(R)$  for ensembles of linear codes and systematic linear codes. Prior to that, we first define these ensembles. A linear code is defined by mapping each of the  $M = 2^K$  binary information (row) vectors  $\mathbf{u}_m$ ,  $0 \leq m \leq M - 1$ , of length  $K$ , into its corresponding code (row) vector  $\mathbf{v}_m$ , of length  $N$ , in the following way:

$$\mathbf{v}_m = \mathbf{u}_m \mathbf{G} \oplus \mathbf{v}_0, \quad m = 0, 1, \dots, M - 1,$$

where  $\mathbf{G}$  is a binary generator matrix of dimension  $K \times N$  and  $\mathbf{v}_0$  is an additive vector of length  $N$ . The  $\oplus$  operation denotes a summation modulo 2 and the multiplication between  $\mathbf{u}_m$  and  $\mathbf{G}$  is conducted over the field  $GF(2)$ . A systematic linear code is defined in the same manner, with the restriction that the left  $K \times K$  block of  $\mathbf{G}$  (the systematic part of  $\mathbf{G}$ ) forms the identity matrix (thus, the first  $K$  bits of each code vector,  $\mathbf{v}_m$ , form the corresponding information vector,  $\mathbf{u}_m$ ).

We now consider a random coding distribution, which is i.i.d. over the ensemble of linear codes (or systematic linear codes), for which the elements of  $\mathbf{G}$  (or  $\tilde{\mathbf{G}}$ , the non-systematic part of  $\mathbf{G}$ , in the case of systematic linear codes) and  $\mathbf{v}_0$  are drawn independently using a uniform single-letter distribution  $Q^* = \left\{\frac{1}{2}, \frac{1}{2}\right\}$  (fair coin tossing). We also define the family of the *binary-input, output-symmetric (BIOS) channels*, as channels with a binary input alphabet  $\mathcal{X}$  ("0" and "1"), an output alphabet  $\mathcal{Y}$  (possibly infinite), where the transition probabilities satisfy  $P(y|0) = P(-y|1), \forall y \in \mathcal{Y}$ , for a well defined operation "-" (note that the definition of symmetry can be used as long as each  $y \in \mathcal{Y}$  satisfies that  $-y \in \mathcal{Y}$  as well). For example, the BSC, when mapping "0"  $\rightarrow +1$  and "1"  $\rightarrow -1$ , is a BIOS channel. The additive Gaussian channel with two antipodal input letters,  $x_1$  and  $x_2$ , is also a BIOS channel.

The following theorem is stated with regard to codes governed by the above mentioned ensembles and transmitted via a BIOS channel:

**Theorem 2** *Consider the sequence of ensembles of linear or systematic linear codes, where the elements of  $\mathbf{G}$  (or  $\tilde{\mathbf{G}}$ ) and  $\mathbf{v}_0$  are drawn independently by fair coin tossing. Let  $\{P_\theta, \theta \in \Theta\}$  be a family of BIOS DMC's. Then, the lower bound  $\xi_{LB}^*(R)$  of eq. (12), continues to hold, with  $\Delta^*(P) = \ln 2 - H(P)$ .*

Theorem 2 is proved in Section A.1 of the Appendix.

The single-letter expression derivation for  $\xi_{LB}^*(R)$  is carried out (see Section A.1 of the Appendix) using the same techniques as in Gallager's classical work, which are tight in the random coding sense. We therefore believe that the achievable lower bounds to the real exponential rates are tight as well. To demonstrate the tightness of the lower bounds suggested in (12) (for general codes) and in Theorem 2 (for linear and systematic linear codes), we have the following lemma:

**Lemma 1** *Consider the family of BSC's parameterized by the crossover probability  $\theta$ . Then,  $\xi_{LB}^*(R) = 1$  and hence  $\xi^*(R) = 1$ , in the following cases:*

- (i) *The ensemble of all codes with  $Q_N(\mathbf{x}) = 2^{-N}$  for all  $\mathbf{x}$ .*
- (ii) *The ensemble of linear codes and systematic linear codes, as in Theorem 2, with  $\Delta^*(P) = \ln 2 - H(P)$ .*

Lemma 1 is proved in Section V.

It should be mentioned that proving that under the BSC model  $\xi^* = 1$  is universally achievable by random coding over general codes and linear codes is by no means new, as it was already proved and discussed in [1]. Nevertheless, it demonstrates the tightness of  $\xi_{LB}^*(R)$ . However, to the best of our knowledge, the same result regarding ensembles of systematic linear codes has not been proved yet and is first shown here.

### 3.3 Convolutional codes

For the special case of the BSC mentioned above, we now introduce the following result, related to ensembles of time-varying convolutional codes, when the minimax decoding is used. Prior to that, we first define this ensemble and the bit error exponent related to it.

A convolutional code of rate  $b/n$  ( $b, n$  – positive integers) and constraint length  $Kb$  is defined as one for which at each time instant  $t \geq 0$ , the code vector of length  $n$ ,  $\mathbf{v}_t$ , is obtained by

$$\mathbf{v}_t = \sum_{j=0}^{\min\{t, Kb-1\}} \mathbf{u}_{t-j} \mathbf{G}_j \oplus \mathbf{v}_0, \quad (13)$$

where  $\mathbf{u}_{t-j}$  is a binary information row vector of length  $b$  at time  $t-j$ ,  $\mathbf{G}_j, 0 \leq j \leq K-1$ , are binary matrices with  $b$  rows and  $n$  columns each, and  $\mathbf{v}_0$  is a vector of length  $n$ .

Let us now consider a code  $\mathcal{C}$ , governed by i.i.d. random coding over the ensemble of time-varying convolutional codes, whose code vector of time instant  $t \geq 0$ ,  $\mathbf{v}_t$ , is obtained by

$$\mathbf{v}_t = \sum_{j=0}^{\min\{t, K-1\}} \mathbf{u}_{t-j} \mathbf{G}_j^t \oplus \mathbf{v}_0^t, \quad (14)$$

where at each time instant  $t$ , the elements of  $\mathbf{G}_j^t, 0 \leq j \leq K-1$  and  $\mathbf{v}_0^t$  are drawn independently using the uniform single-letter distribution  $\{\frac{1}{2}, \frac{1}{2}\}$ .

The average bit error probability,  $\overline{P_b(\Omega_K)}$ , associated with a sequence of decoders  $\Omega_K = \{\Omega_{K,N}\}_{N=1}^{\infty}$  of block length  $N$  and constraint length  $K$ , and averaged over the ensemble of time-varying convolutional codes, is defined as the expected relative frequency of bit errors in the decoded information stream, i.e.

$$\overline{P_b(\Omega_K)} = \limsup_{N \rightarrow \infty} \overline{P_b(\Omega_{K,N})}. \quad (15)$$

The bit error exponent associated with a sequence of decoders  $\Omega = \{\Omega_K\}_{K=1}^{\infty}$  is defined as

$$\overline{E_b(\Omega)} = - \limsup_{K \rightarrow \infty} \frac{1}{K} \ln \overline{P_b(\Omega_K)}. \quad (16)$$

**Theorem 3** *Consider the sequence of ensembles of time-varying convolutional codes of rate  $b/n$  and constraint length  $Kb$  (with  $K \rightarrow \infty$ ), described as in the previous paragraph, and assume a family of BSC's parameterized by the crossover probability  $\theta$ .*

*The achievable bit error exponent (as defined in (16)) using the minimax decoder is equal to the one when  $\theta$  is known and the ML decoder is used.*

The proof of this theorem is based on the following observation:

Under the BSC model with an unknown crossover probability  $\theta$ , the minimax decision rule (as defined in (3)) is equivalent to a decision rule, denoted by  $\Lambda$ , and defined as:

$$\Lambda_m = \{\mathbf{y} | \rho(\mathbf{x}_m, \mathbf{y}) \leq \rho(\mathbf{x}_{m'}, \mathbf{y}), \forall m' \neq m\}, \quad (17)$$

with ties broken arbitrarily, where

$$\rho(\mathbf{x}, \mathbf{y}) = \min \{\delta(\mathbf{x}, \mathbf{y}), 1 - \delta(\mathbf{x}, \mathbf{y})\}. \quad (18)$$

As mentioned in Section II,  $\delta(\mathbf{x}, \mathbf{y})$  denotes the normalized Hamming distance between  $\mathbf{x}$  and  $\mathbf{y}$ . This equivalence is proved in Section A.7 of the Appendix. We should note that

for this case, the minimax decoder coincides with the MMI decoder as well. Based on this equivalence, the full proof of Theorem 3 is given in Section VII. We also introduce an efficient implementation of minimax decoding, based on a slightly modified version of the Viterbi algorithm. This is done by applying the Viterbi algorithm twice: first for minimum Hamming distance, and then for maximum Hamming distance. This process results in two survivors and the selection between them is done in favor of the one whose normalized Hamming metric is more distant from  $\frac{1}{2}$  (the one with the minimal  $\rho$ ).

## 4 Proof of Theorem 1

We first observe that for a DMC,  $\{P_\theta(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$ , and for each vector pair  $(\mathbf{x}, \mathbf{y})$ , the minimax metric for a given  $\theta$ ,  $f_\theta(\mathbf{x}, \mathbf{y})$ , depends on  $\mathbf{x}$  and  $\mathbf{y}$  only via their joint empirical distribution:

$$f_\theta(\mathbf{x}, \mathbf{y}) = \hat{E}\mathbf{x}\mathbf{y} \ln P_\theta(Y|X) + \xi E_r^*(\theta). \quad (19)$$

We, therefore, conclude that the value of  $\theta$  maximizing  $f_\theta(\mathbf{x}, \mathbf{y})$  also depends on  $\mathbf{x}$  and  $\mathbf{y}$  only via their joint empirical distribution. Let  $\Theta_N$  denote the subset of  $\Theta$  with values of  $\theta$  that achieve  $\max_\theta f_\theta(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}, \mathbf{y})$  as  $(\mathbf{x}, \mathbf{y})$  exhaust  $\mathcal{X}^N \times \mathcal{Y}^N$ . In the decoding process, maximization over  $\theta$  can be achieved only by points in  $\Theta_N$ . Since the number of joint empirical distributions of  $(\mathbf{x}, \mathbf{y})$  is upper bounded by  $(N+1)^{|\mathcal{X}||\mathcal{Y}|}$ , then  $|\Theta_N| \leq (N+1)^{|\mathcal{X}||\mathcal{Y}|}$  as well.

As a first step, we assume given channel input and output vectors,  $\mathbf{x}$  and  $\mathbf{y}$ , respectively. Considering a random coding distribution,  $Q_N$ , we exponentially evaluate the probability of having another codeword  $\mathbf{x}'$  that is preferred by the minimax decoder over  $\mathbf{x}$ . This probability will be denoted by  $a(\mathbf{x}, \mathbf{y})$ .

$$\begin{aligned} a(\mathbf{x}, \mathbf{y}) &= Q_N \{f(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y})\} \\ &= Q_N \left\{ \max_{\theta' \in \Theta_N} f_{\theta'}(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y}) \right\} \\ &\stackrel{(a)}{=} \max_{\theta' \in \Theta_N} Q_N \{f_{\theta'}(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y})\} \\ &= \max_{\theta' \in \Theta_N} Q_N \left\{ \sum_{i=1}^N \ln P_{\theta'}(y_i | X'_i) \geq -N\xi E_r^*(\theta') + N \cdot f(\mathbf{x}, \mathbf{y}) \right\} \\ &\stackrel{(b)}{=} \max_{\theta' \in \Theta_N} \min_{s \geq 0} E_{Q_N} \left[ \exp \left\{ s \left[ \sum_{i=1}^N \ln P_{\theta'}(y_i | X'_i) + N\xi E_r^*(\theta') \right] \right\} \right] \cdot \exp \{-sNf(\mathbf{x}, \mathbf{y})\} \end{aligned}$$

$$= \max_{\theta' \in \Theta_N} \min_{s \geq 0} E_{Q_N} e^{sN f_{\theta'}(\mathbf{X}', \mathbf{y})} \cdot e^{-sN f(\mathbf{x}, \mathbf{y})}, \quad (20)$$

where (a) is true since

$$\begin{aligned} \max_{\theta' \in \Theta_N} Q_N \{f_{\theta'}(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y})\} &\leq Q_N \left\{ \max_{\theta' \in \Theta_N} f_{\theta'}(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y}) \right\} \\ &= Q_N \left\{ \bigcup_{\theta' \in \Theta_N} f_{\theta'}(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y}) \right\} \\ &\leq \sum_{\theta' \in \Theta_N} Q_N \{f_{\theta'}(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y})\} \\ &= |\Theta_N| \cdot \max_{\theta' \in \Theta_N} Q_N \{f_{\theta'}(\mathbf{X}', \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y})\}, \quad (21) \end{aligned}$$

and in (b) we used the Cherooff bound, which is tight in the exponential sense.

By using the method of types, it is proved in Section A.3 of the Appendix that for any real  $\alpha$ ,

$$E_{Q_N}[e^{N\alpha f_{\theta}(\mathbf{X}, \mathbf{y})}] = e^{N[\alpha \xi E_r^*(\theta) - \min_P \mathbf{x} | \mathbf{y}^A(\theta, \alpha, P \mathbf{x} \mathbf{y})]}, \quad (22)$$

where the function  $A(\theta, \alpha, P_{xy})$  is defined as in (8).

Using this observation, we can continue to evaluate  $a(\mathbf{x}, \mathbf{y})$  as follows:

$$\begin{aligned} a(\mathbf{x}, \mathbf{y}) &\doteq \max_{\theta' \in \Theta_N} \min_{s \geq 0} \exp \left\{ N[s \xi E_r^*(\theta') - \min_{\mathbf{x}' | \mathbf{y}} A(\theta', s, P_{\mathbf{x}' \mathbf{y}})] \right\} \cdot \exp \{-sN f(\mathbf{x}, \mathbf{y})\} \\ &= \max_{\theta' \in \Theta_N} \min_{s \geq 0} \exp \left\{ -N[-s \xi E_r^*(\theta') + \min_{\mathbf{x}' | \mathbf{y}} A(\theta', s, P_{\mathbf{x}' \mathbf{y}}) + s f(\mathbf{x}, \mathbf{y})] \right\} \\ &\triangleq \max_{\theta' \in \Theta_N} \min_{s \geq 0} \exp \left\{ -N[G(\theta', s, \xi, P \mathbf{x} \mathbf{y})] \right\}. \quad (23) \end{aligned}$$

Therefore, the probability that the decoder will prefer any of the other  $M - 1$  codevectors rather than the transmitted codevector  $\mathbf{x}$  can be evaluated as follows:

$$\begin{aligned} 1 - (1 - a(\mathbf{x}, \mathbf{y}))^{M-1} &\stackrel{(a)}{=} \min \{1, e^{NR} \cdot a(\mathbf{x}, \mathbf{y})\} \\ &= \min \left\{ 1, \max_{\theta' \in \Theta_N} \min_{s \geq 0} \exp \{-N[G(\theta', s, \xi, P \mathbf{x} \mathbf{y}) - R]\} \right\} \\ &= \max_{\theta' \in \Theta_N} \min_{s \geq 0} \exp \left\{ -N \cdot \max_{0 \leq \rho \leq 1} \rho [G(\theta', s, \xi, P \mathbf{x} \mathbf{y}) - R] \right\} \\ &= \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \exp \left\{ -N \cdot [\rho G(\theta', s, \xi, P \mathbf{x} \mathbf{y}) - \rho R] \right\}, \quad (24) \end{aligned}$$

where the equivalence in (a) (see [9], Section V, and [8], Section A.2 p. 109-110) implies that the union bound in the random coding error exponent is tight.

Now, we will evaluate  $\overline{S}_N$ , the average of the minimax criterion over the ensemble of codebooks governed by a random coding distribution, for the minimax decoder defined in (3):

$$\begin{aligned}
\overline{S}_N &\doteq \max_{\theta \in \Theta} \left\{ \frac{\overline{P}_E(\Omega|\theta)}{e^{-N\xi E_r^*(\theta)}} \right\} \\
&= \max_{\theta \in \Theta} \left\{ e^{N\xi E_r^*(\theta)} \sum_{\mathbf{x} \in \mathcal{X}^N} Q_N(\mathbf{x}) \sum_{\mathbf{y} \in \mathcal{Y}^N} P_\theta(\mathbf{y}|\mathbf{x}) \left[ 1 - (1 - a(\mathbf{x}, \mathbf{y}))^{M-1} \right] \right\} \\
&\doteq \max_{\theta \in \Theta} \left\{ \sum_{\mathbf{x} \in \mathcal{X}^N} Q_N(\mathbf{x}) \sum_{\mathbf{y} \in \mathcal{Y}^N} e^{Nf_\theta(\mathbf{x}, \mathbf{y})} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \exp \left\{ -N[\rho G(\theta', s, \xi, P_{\mathbf{x}\mathbf{y}}) - \rho R] \right\} \right\} \\
&\stackrel{(a)}{=} \max_{\theta \in \Theta} \left\{ \sum_{T_{\mathbf{x}\mathbf{y}} \subset \mathcal{X}^N \times \mathcal{Y}^N} Q_N(T_{\mathbf{x}}) |T_{\mathbf{y}}| e^{Nf_\theta(\mathbf{x}, \mathbf{y})} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \right. \\
&\quad \left. e^{N\rho s \xi E_r^*(\theta')} \cdot e^{-N\rho \min_P} \mathbf{x}'|_{\mathbf{y}} A(\theta', s, P_{\mathbf{x}'\mathbf{y}}) \cdot e^{-N\rho s f(\mathbf{x}, \mathbf{y})} \cdot e^{N\rho R} \right\} \\
&\doteq \max_{P_{\mathbf{x}\mathbf{y}}} \left\{ e^{-N\Delta_N^*(P_{\mathbf{x}})} \cdot e^{NH\mathbf{x}\mathbf{y}(Y|X)} \left[ \max_{\theta \in \Theta_N} e^{Nf_\theta(\mathbf{x}, \mathbf{y})} \right] \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \max_{P_{\mathbf{x}'|_{\mathbf{y}}}} \right. \\
&\quad \left. e^{N\rho s \xi E_r^*(\theta')} \cdot e^{-N\rho A(\theta', s, P_{\mathbf{x}'\mathbf{y}})} \cdot \left[ \max_{\theta'' \in \Theta_N} e^{Nf_{\theta''}(\mathbf{x}, \mathbf{y})} \right]^{-\rho s} \cdot e^{N\rho R} \right\} \\
&\stackrel{(b)}{=} \max_{P_{\mathbf{x}\mathbf{y}}} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \max_{P_{\mathbf{x}'|_{\mathbf{y}}}} \left\{ e^{-N\Delta^*(P_{\mathbf{x}})} e^{NH\mathbf{x}\mathbf{y}(Y|X)} e^{N\rho s \xi E_r^*(\theta')} e^{-N\rho A(\theta', s, P_{\mathbf{x}'\mathbf{y}})} \right. \\
&\quad \left. \left[ \max_{\theta \in \Theta_N} e^{Nf_\theta(\mathbf{x}, \mathbf{y})} \right]^{1-\rho s} e^{N\rho R} \right\} \\
&= \max_{P_{\mathbf{x}\mathbf{y}}} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \max_{P_{\mathbf{x}'|_{\mathbf{y}}}} \left\{ \exp \left\{ N[-\Delta^*(P_{\mathbf{x}}) + H\mathbf{x}\mathbf{y}(Y|X) + \rho s \xi E_r^*(\theta') \right. \right. \\
&\quad \left. \left. - \rho A(\theta', s, P_{\mathbf{x}'\mathbf{y}}) + \rho R] \right\} \left[ \max_{\theta \in \Theta_N} e^{Nf_\theta(\mathbf{x}, \mathbf{y})} \right]^{1-\rho s} \right\} \\
&\triangleq \max_{P_{\mathbf{x}\mathbf{y}}} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \max_{P_{\mathbf{x}'|_{\mathbf{y}}}} \left\{ \exp \left\{ N \cdot T(\theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho, \xi, R) \right\} \right. \\
&\quad \left. \left[ \max_{\theta \in \Theta_N} e^{Nf_\theta(\mathbf{x}, \mathbf{y})} \right]^{1-\rho s} \right\}, \tag{25}
\end{aligned}$$

where in (a) we switched to a summation over the joint empirical types of  $\mathbf{x}$  and  $\mathbf{y}$  (which is legitimate since both  $f_\theta(\mathbf{x}, \mathbf{y})$  and  $G(\theta', s, \xi, P_{\mathbf{x}\mathbf{y}})$  depend on  $\mathbf{x}$  and  $\mathbf{y}$  via their joint empirical distribution), and in (b), we used the convergence assumption of the random coding distributions within the class  $\mathcal{Q}$  to claim that  $\Delta_N^*(P_{\mathbf{x}}) \rightarrow \Delta^*(P_{\mathbf{x}})$  as  $N \rightarrow \infty$  independently of  $P_{\mathbf{x}}$ , and also united the optimizations over  $\theta$  and  $\theta''$ .

We should observe that:

$$\begin{aligned}
& \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 0}} \max_{P \mathbf{x}' | \mathbf{y}} \left\{ \exp \left\{ N \cdot T(\theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) \right\} \left[ \max_{\theta \in \Theta_N} e^{N f_\theta(\mathbf{x}, \mathbf{y})} \right]^{1-\rho s} \right\} = \\
& = \min \left\{ \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \max_{\theta \in \Theta_N} \exp \left\{ N \left[ T(\theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) + f_\theta(\mathbf{x}, \mathbf{y})(1 - \rho s) \right] \right\}, \right. \\
& \quad \left. \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \min_{\theta \in \Theta_N} \exp \left\{ N \left[ T(\theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) + f_\theta(\mathbf{x}, \mathbf{y})(1 - \rho s) \right] \right\} \right\} \\
& \stackrel{(a)}{=} \min \left\{ \max_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \exp \left\{ N \left[ T(\theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) + f_\theta(\mathbf{x}, \mathbf{y})(1 - \rho s) \right] \right\}, \right. \\
& \quad \left. \min_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \exp \left\{ N \left[ T(\theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) + f_\theta(\mathbf{x}, \mathbf{y})(1 - \rho s) \right] \right\} \right\} \tag{26}
\end{aligned}$$

$$\begin{aligned}
& \triangleq \min \left\{ \max_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \exp \left\{ N \cdot \tilde{T}(\theta, \theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) \right\}, \right. \\
& \quad \left. \min_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \exp \left\{ N \cdot \tilde{T}(\theta, \theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) \right\} \right\}, \tag{27}
\end{aligned}$$

where in (a), two interchanges are made: one between the minimization over  $\rho$  and  $s$  and the maximization over  $\theta$  in the left term of the outer minimization, and one between the maximization over  $P \mathbf{x}' | \mathbf{y}$  and the minimization over  $\theta$  in the right term of the outer minimization. The first interchange is justified in the Appendix, Section A.2. The second interchange is possible since the term to be optimized is a product of two exponential terms, one depends on  $P \mathbf{x}' | \mathbf{y}$  and one depends on  $\theta$ , therefore the optimizations can be done independently.

Consequently, we conclude that:

$$\begin{aligned}
\overline{S}_N & \doteq \max_{P \mathbf{x} \mathbf{y}} \max_{\theta' \in \Theta_N} \min \left\{ \max_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \exp \left\{ N \cdot \tilde{T}(\theta, \theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) \right\}, \right. \\
& \quad \left. \min_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \exp \left\{ N \cdot \tilde{T}(\theta, \theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) \right\} \right\} \\
& = \max_{P \mathbf{x} \mathbf{y}} \max_{\theta' \in \Theta_N} \min \left\{ \exp \left\{ N \cdot \max_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \tilde{T}(\theta, \theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) \right\}, \right. \\
& \quad \left. \exp \left\{ N \cdot \min_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P \mathbf{x}' | \mathbf{y}} \tilde{T}(\theta, \theta', P \mathbf{y}, P \mathbf{x} | \mathbf{y}, P \mathbf{x}' | \mathbf{y}, s, \rho, \xi, R) \right\} \right\}
\end{aligned}$$



$$\begin{aligned}
= & \exp \left\{ N \cdot \max_{P_{\mathbf{x}|\mathbf{y}}} \max_{\theta' \in \Theta_N} \min \left\{ \max_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P_{\mathbf{x}'|\mathbf{y}}} \tilde{T}(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho, \xi, R), \right. \right. \\
& \left. \left. \min_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P_{\mathbf{x}'|\mathbf{y}}} \tilde{T}(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho, \xi, R) \right\} \right\}.
\end{aligned} \tag{28}$$

Now,

$$\begin{aligned}
\tilde{T}(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho, \xi, R) &= \\
&= -\Delta^*(P_{\mathbf{x}}) + H_{\mathbf{y}}(Y) - I_{\mathbf{x}\mathbf{y}}(X; Y) + \rho s \xi E_r^*(\theta') - \rho A(\theta', s, P_{\mathbf{x}'|\mathbf{y}}) \\
&\quad + (1 - \rho s) \hat{E}_{\mathbf{x}\mathbf{y}} \ln P_{\theta}(Y|X) + (1 - \rho s) \xi E_r^*(\theta) + \rho R \\
&= -A(\theta, 1 - \rho s, P_{\mathbf{x}\mathbf{y}}) - \rho A(\theta', s, P_{\mathbf{x}'|\mathbf{y}}) + H_{\mathbf{y}}(Y) + \rho s \xi E_r^*(\theta') \\
&\quad + (1 - \rho s) \xi E_r^*(\theta) + \rho R \\
&= -B(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho) + \rho s \xi E_r^*(\theta') + (1 - \rho s) \xi E_r^*(\theta) + \rho R,
\end{aligned} \tag{29}$$

where the function  $B(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho)$  is defined as in (9).

Therefore, in order for  $\overline{S}_N$  to grow sub-exponentially with  $N$ , we seek the maximal  $\xi$  such that:

$$\begin{aligned}
\max_{P_{\mathbf{x}|\mathbf{y}}} \max_{\theta' \in \Theta_N} \min \left\{ \max_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P_{\mathbf{x}'|\mathbf{y}}} \tilde{T}(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho, \xi, R), \right. \\
\left. \min_{\theta \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P_{\mathbf{x}'|\mathbf{y}}} \tilde{T}(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho, \xi, R) \right\} \leq 0 \tag{30}
\end{aligned}$$

As the empirical distributions become dense in continuum of probability distributions as  $N \rightarrow \infty$ , and since the function  $\tilde{T}(\theta, \theta', P_{\mathbf{y}}, P_{\mathbf{x}|\mathbf{y}}, P_{\mathbf{x}'|\mathbf{y}}, s, \rho, \xi, R)$  is continuous in  $P_{\mathbf{y}}$ ,  $P_{\mathbf{x}|\mathbf{y}}$  and  $P_{\mathbf{x}'|\mathbf{y}}$ , it is equivalent to perform the above optimizations over continuous distributions rather than empirical distributions. The same token can be used in order to broaden the maximization space for  $\theta$  and  $\theta'$  from  $\Theta_N$  to  $\Theta$ . Thus, the condition becomes:

$$\begin{aligned}
\max_{P_{X|Y}} \max_{\theta' \in \Theta} \min \left\{ \max_{\theta \in \Theta} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P_{X'|Y}} \tilde{T}(\theta, \theta', P_y, P_{X|Y}, P_{X'|Y}, s, \rho, \xi, R), \right. \\
\left. \min_{\theta \in \Theta} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P_{X'|Y}} \tilde{T}(\theta, \theta', P_y, P_{X|Y}, P_{X'|Y}, s, \rho, \xi, R) \right\} \leq 0 \tag{31}
\end{aligned}$$

In other words, a maximal  $\xi$  is sought such that:

$$\forall P_{Xy}, \forall \theta' \in \Theta$$

$$\max_{\theta \in \Theta} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \max_{P_{X'|y}} \tilde{T}(\theta, \theta', P_y, P_{X|y}, P_{X'|y}, s, \rho, \xi, R) \leq 0 \quad (32)$$

or

$$\min_{\theta \in \Theta} \min_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \max_{P_{X'|y}} \tilde{T}(\theta, \theta', P_y, P_{X|y}, P_{X'|y}, s, \rho, \xi, R) \leq 0 \quad (33)$$

An equivalent condition is:

$$\forall P_{Xy}, \forall \theta' \in \Theta$$

$$\xi \leq \min_{\theta \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \min_{P_{X'|y}} \frac{B(\theta, \theta', P_y, P_{X|y}, P_{X'|y}, s, \rho) - \rho R}{(1 - \rho s)E_r^*(\theta) + \rho s E_r^*(\theta')} \quad (34)$$

or

$$\xi \leq \max_{\theta \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \min_{P_{X'|y}} \frac{B(\theta, \theta', P_y, P_{X|y}, P_{X'|y}, s, \rho) - \rho R}{(1 - \rho s)E_r^*(\theta) + \rho s E_r^*(\theta')} \quad (35)$$

Therefore,

$$\begin{aligned} \xi^*(R) = & \min_{P_{Xy}} \min_{\theta' \in \Theta} \max \left\{ \min_{\theta \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq s \leq 1/\rho}} \min_{P_{X'|y}} \frac{B(\theta, \theta', P_y, P_{X|y}, P_{X'|y}, s, \rho) - \rho R}{(1 - \rho s)E_r^*(\theta) + \rho s E_r^*(\theta')}, \right. \\ & \left. \max_{\theta \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ s \geq 1/\rho}} \min_{P_{X'|y}} \frac{B(\theta, \theta', P_y, P_{X|y}, P_{X'|y}, s, \rho) - \rho R}{(1 - \rho s)E_r^*(\theta) + \rho s E_r^*(\theta')} \right\}. \quad (36) \end{aligned}$$

## 5 Example - the BSC

In this section, we demonstrate that for the special case of BSC with an unknown crossover probability, and a uniform random coding distribution,  $\xi_{LB}^*(R) = 1$  and hence  $\xi^*(R) = 1$ , in agreement with well known results [1].

Consider the lower bound (12) and choose the uniform single-letter random coding distribution  $Q^* = \{\frac{1}{2}, \frac{1}{2}\}$ .

Now, the value of  $A(\theta, \alpha, P_{XY})$  is (see (11)):

$$A(\theta, \alpha, P_{XY}) = \ln 2 - H(X|Y) - \alpha E \ln P_\theta(Y|X) \quad (37)$$

Therefore,

$$\min_{P_{X|Y}} A(\theta, \alpha, P_{XY}) = \ln 2 - \max_{P_{X|Y}} \{H(X|Y) + \alpha E \ln P_\theta(Y|X)\} \quad (38)$$

In addition, for the case of BSC with an unknown crossover probability,  $\theta$ , we have (see [7], Section VI):

$$\begin{aligned} \max_{P_{X|Y}} \{H(X|Y) + \alpha E \ln P_\theta(Y|X)\} &= \ln[(1 - \theta)^\alpha + \theta^\alpha] \\ &\triangleq \mathcal{V}(\theta, \alpha) \end{aligned} \quad (39)$$

From these two observations, we conclude that:

$$\min_{P_{X|Y}} A(\theta, \alpha, P_{XY}) = \ln 2 - \mathcal{V}(\theta, \alpha) \quad (40)$$

Using (9), we get:

$$\begin{aligned} \xi_{LB}^*(R) &= \min_{P_Y} \min_{\theta, \theta' \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \min_{P_{X|Y}} \min_{P_{X'|Y}} \frac{A(\theta, 1 - \lambda\rho, P_{XY}) + \rho A(\theta', \lambda, P_{X'Y}) - H(Y) - \rho R}{(1 - \lambda\rho) \cdot E_r^*(\theta) + \lambda\rho \cdot E_r^*(\theta')} \\ &= \min_{P_Y} \min_{\theta, \theta' \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \frac{(1 + \rho) \ln 2 - \mathcal{V}(\theta, 1 - \lambda\rho) - \rho \mathcal{V}(\theta', \lambda) - H(Y) - \rho R}{(1 - \lambda\rho) \cdot E_r^*(\theta) + \lambda\rho \cdot E_r^*(\theta')} \\ &\geq \min_{\theta, \theta' \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \min_{P_Y} \frac{(1 + \rho) \ln 2 - \mathcal{V}(\theta, 1 - \lambda\rho) - \rho \mathcal{V}(\theta', \lambda) - H(Y) - \rho R}{(1 - \lambda\rho) \cdot E_r^*(\theta) + \lambda\rho \cdot E_r^*(\theta')} \\ &= \min_{\theta, \theta' \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \frac{\rho \ln 2 - \mathcal{V}(\theta, 1 - \lambda\rho) - \rho \mathcal{V}(\theta', \lambda) - \rho R}{(1 - \lambda\rho) \cdot E_r^*(\theta) + \lambda\rho \cdot E_r^*(\theta')}. \end{aligned} \quad (41)$$

Now, the random coding error exponent associated with ML decoding,  $E_r^*(\theta)$ , to which the minimax decoding error exponent is compared, is achieved for the BSC model by the following optimization (see [10, Sect. 3.1, 3.2 and 3.4]):

$$\begin{aligned} E_r^*(\theta) &= \max_{0 \leq \rho \leq 1} \max_Q \left\{ -\ln \sum_{y \in \{0,1\}} \left[ \sum_{x \in \{0,1\}} Q(x) \cdot P_\theta(y|x)^{\frac{1}{1+\rho}} \right]^{1+\rho} - \rho R \right\} \\ &\stackrel{(a)}{=} \max_{0 \leq \rho \leq 1} \left\{ \rho \ln 2 - (1 + \rho) \ln \left[ (1 - \theta)^{\frac{1}{1+\rho}} + \theta^{\frac{1}{1+\rho}} \right] - \rho R \right\} \\ &= \max_{0 \leq \rho \leq 1} \left\{ \rho \ln 2 - (1 + \rho) \mathcal{V} \left( \theta, \frac{1}{1 + \rho} \right) - \rho R \right\} \\ &\triangleq \max_{0 \leq \rho \leq 1} E_r(\theta, \rho), \end{aligned} \quad (42)$$

where in (a), the inner maximization is achieved by taking  $Q^* = \{\frac{1}{2}, \frac{1}{2}\}$  ([10, Sect. 3.4]).

Let us now define  $\rho' = \frac{\lambda\rho}{1 - \lambda\rho}$  and  $\rho'' = \frac{1}{\lambda} - 1$ , and rewrite the numerator of (41) as follows:

$$\rho \ln 2 - \mathcal{V}(\theta, 1 - \lambda\rho) - \rho \mathcal{V}(\theta', \lambda) - \rho R =$$

$$\begin{aligned}
&= \rho \ln 2 - \mathcal{V}\left(\theta, \frac{1}{1 + \rho'}\right) - \rho \mathcal{V}\left(\theta', \frac{1}{1 + \rho''}\right) - \rho R \\
&= (1 - \lambda\rho) \left[ \rho' \ln 2 - (1 + \rho') \mathcal{V}\left(\theta, \frac{1}{1 + \rho'}\right) - \rho' R \right] + \\
&\quad \lambda\rho \left[ \rho'' \ln 2 - (1 + \rho'') \mathcal{V}\left(\theta', \frac{1}{1 + \rho''}\right) - \rho'' R \right] \\
&= (1 - \lambda\rho) \cdot E_r(\theta, \rho') + \lambda\rho \cdot E_r(\theta', \rho'') \\
&= (1 - \lambda\rho) \cdot E_r\left(\theta, \frac{\lambda\rho}{1 - \lambda\rho}\right) + \lambda\rho \cdot E_r\left(\theta', \frac{1}{\lambda} - 1\right). \tag{43}
\end{aligned}$$

Finally, we get that

$$\xi_{LB}^*(R) \geq \min_{\theta, \theta' \in \Theta} \max_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \frac{(1 - \lambda\rho) \cdot E_r(\theta, \frac{\lambda\rho}{1 - \lambda\rho}) + \lambda\rho \cdot E_r(\theta', \frac{1}{\lambda} - 1)}{(1 - \lambda\rho) \cdot E_r^*(\theta) + \lambda\rho \cdot E_r^*(\theta')}. \tag{44}$$

Now, by choosing  $\lambda = \frac{1}{1 + \tilde{\rho}}$ , where  $\tilde{\rho}$  is the achiever of  $E_r^*(\theta') = \max_{0 \leq \rho \leq 1} E_r(\theta', \rho)$ , and  $\rho = \frac{\hat{\rho}}{1 + \hat{\rho}}(1 + \tilde{\rho})$ , where  $\hat{\rho}$  is the achiever of  $E_r^*(\theta) = \max_{0 \leq \rho \leq 1} E_r(\theta, \rho)$  (observing that  $\frac{\hat{\rho}}{1 + \hat{\rho}}(1 + \tilde{\rho}) \leq 1$ , therefore this choice is feasible), we get that both the numerator and the denominator of (44) equal to  $(1 - \lambda\rho) \cdot E_r(\theta, \hat{\rho}) + \lambda\rho \cdot E_r(\theta', \tilde{\rho})$ , and so,  $\xi_{LB}^*(R) = 1$ .

We should note that for the BSC model, the same conclusion (i.e.,  $\xi^* = 1$ ) holds also for linear codes and systematic linear codes (as the optimal random coding distribution that was used is  $Q^* = \{\frac{1}{2}, \frac{1}{2}\}$  (see (42)).

## 6 Proof of Theorem 3

First, consider a given channel output related to the entire transmitted sequence of information. Without loss of generality, the all-zero message will be assumed to be transmitted. Let us now consider a segment of length  $K + l$ ,  $l \geq 0$ , of the transmitted information vector, and any other incorrect path diverging from it at node  $j$  and emerging at node  $j + K + l$  (note that the minimum length of a diverging path is  $K$  since after a non-zero vector is inserted to the encoder,  $K - 1$  zero vectors are needed in order to return to the all-zero state).

We observe that the information sequence related to such an incorrect path has the following structure (we ignore the values of the information sequence outside the range

$(j, j + K + l - 1))$ :

$$\mathbf{u}_j, \mathbf{u}_{j+1}, \dots, \mathbf{u}_{j+l}, \underbrace{\mathbf{0}, \dots, \mathbf{0}}_{K-1}$$

where all of the vectors are of length  $b$ .

In order for the incorrect path to diverge exactly from node  $j$  to node  $j + K + l$ ,  $\mathbf{u}_j$  and  $\mathbf{u}_{j+l}$  can be any of the  $2^b - 1$  non-zero vectors (thus, there are  $(2^b - 1)^2$  possibilities for their values), and each of the  $l - 1$  information vectors  $\mathbf{u}_{j+1}, \dots, \mathbf{u}_{j+l-1}$  can be any binary vector of length  $b$ , with the restriction of no more than  $K - 2$  consecutive all-zero vectors (thus, there are less than  $2^{b(l-1)}$  possibilities for their values). Therefore, the number of such incorrect paths, denoted by  $M$ , is upper-bounded by

$$M \leq (2^b - 1)^2 2^{b(l-1)} \leq (2^b - 1) 2^{bl} \quad (45)$$

We next upper bound the probability that an incorrect path is preferred by the minimax decoder (minimizing the metric  $\rho$ ) over the correct path, and then average this probability over the ensemble of time-varying convolutional codes.

We will use  $\mathbf{V}_j = [\mathbf{v}_j, \mathbf{v}_{j+1}, \dots, \mathbf{v}_{j+K+l-1}]$  to denote the code vector of length  $N = n(K+l)$  that corresponds to the correct all-zeros path, while  $\mathbf{V}_j'$  and  $\mathbf{V}_j''$  will be used to denote code vectors that correspond to other incorrect paths. The notation  $\overline{\mathbf{V}}_j$  will be used for the complement vector of  $\mathbf{V}_j$ . A segment of length  $N$  of the corresponding channel output will be denoted by  $\mathbf{Y}_j$ , and  $Q^*$  will be used to denote the random coding distribution.

$$\begin{aligned} & \Pr \left\{ \overline{\rho(\mathbf{V}_j', \mathbf{Y}_j) \leq \rho(\mathbf{V}_j, \mathbf{Y}_j)} \middle| \theta \right\} = \\ &= \sum_{\mathbf{V}_j, \mathbf{V}_j'} Q^*(\mathbf{V}_j, \mathbf{V}_j') \Pr \left\{ \rho(\mathbf{V}_j', \mathbf{Y}_j) \leq \rho(\mathbf{V}_j, \mathbf{Y}_j) \middle| \theta \right\} \\ &\stackrel{(a)}{=} 2^{-2N} \sum_{\mathbf{V}_j, \mathbf{V}_j'} \Pr \left\{ \rho(\mathbf{V}_j', \mathbf{Y}_j) \leq \rho(\mathbf{V}_j, \mathbf{Y}_j) \right\} \\ &= 2^{-2N} \sum_{\mathbf{V}_j, \mathbf{V}_j'} \Pr \left\{ \min \left\{ \delta(\mathbf{V}_j', \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j', \mathbf{Y}_j) \right\} \leq \min \left\{ \delta(\mathbf{V}_j, \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \right\} \\ &= 2^{-2N} \sum_{\mathbf{V}_j, \mathbf{V}_j'} \Pr \left\{ \left[ \delta(\mathbf{V}_j', \mathbf{Y}_j) \leq \min \left\{ \delta(\mathbf{V}_j, \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \right] \right. \\ &\quad \left. \cup \left[ \delta(\overline{\mathbf{V}}_j', \mathbf{Y}_j) \leq \min \left\{ \delta(\mathbf{V}_j, \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \right] \right\} \\ &\stackrel{(b)}{\leq} 2^{-2N} \sum_{\mathbf{V}_j, \mathbf{V}_j'} \Pr \left\{ \delta(\mathbf{V}_j', \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j) \cup \delta(\overline{\mathbf{V}}_j', \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \end{aligned}$$

$$\begin{aligned}
&\stackrel{(c)}{\leq} 2^{-2N} \sum_{\mathbf{V}_j} \sum_{\mathbf{V}'_j} \Pr\{\delta(\mathbf{V}'_j, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j)\} + 2^{-2N} \sum_{\mathbf{V}_j} \sum_{\mathbf{V}'_j} \Pr\{\delta(\overline{\mathbf{V}}'_j, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j)\} \\
&\stackrel{(d)}{=} 2^{-2N} \sum_{\mathbf{V}_j} \sum_{\mathbf{V}'_j} \Pr\{\delta(\mathbf{V}'_j, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j)\} + 2^{-2N} \sum_{\mathbf{V}_j} \sum_{\mathbf{V}''_j} \Pr\{\delta(\mathbf{V}''_j, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j)\} \\
&= 2 \cdot 2^{-2N} \sum_{\mathbf{V}_j} \sum_{\mathbf{V}'_j} \Pr\{\delta(\mathbf{V}'_j, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j)\} \\
&\stackrel{(e)}{\leq} 2 \cdot 2^{-2N} \sum_{\mathbf{V}_j} \sum_{\mathbf{V}'_j} \sum_{\mathbf{Y}_j} \sqrt{P_\theta(\mathbf{Y}_j|\mathbf{V}_j)P_\theta(\mathbf{Y}_j|\mathbf{V}'_j)} \\
&\stackrel{(f)}{=} 2 \left\{ \sum_y \left[ \sum_v \frac{1}{2} \sqrt{P_\theta(y|v)} \right]^2 \right\}^N \\
&\stackrel{=}{=} e^{-NR_{\theta,0}(Q^*)}, \tag{46}
\end{aligned}$$

where

$$R_{\theta,0}(Q^*) = -\ln \sum_y \left[ \sum_v \frac{1}{2} \sqrt{P_\theta(y|v)} \right]^2.$$

In (a) we used the fact that both  $\mathbf{V}_j$  and  $\mathbf{V}'_j$  can attain each of their  $2^N$  possible values equiprobably and independently. This claim for  $\mathbf{V}_j$  (which corresponds to the all-zero path) can be justified due to the fact that the elements of  $\mathbf{G}_j^t, 0 \leq j \leq K-1$  and  $\mathbf{v}_0^t$  are repeatedly randomized at each time instant (see (14)). Therefore,  $\forall 0 \leq i \leq K+l, \mathbf{v}_{j+i} = \mathbf{v}_0^{j+i}$ , thus each one of these vectors is likely to attain each of its  $2^n$  values equiprobably. This claim for  $\mathbf{V}'_j$  (which correspond to the incorrect path) can be justified since  $\mathbf{u}_j$  and  $\mathbf{u}_{j+l}$  are non-zero and  $\mathbf{u}_{j+1}, \dots, \mathbf{u}_{j+l-1}$  cannot include more than  $K-2$  consecutive all-zero vectors. Thus, each code vector of  $\mathbf{v}'_{j+i}, 0 \leq i \leq K+l$  is formed by the modulo-2 sum of  $\mathbf{v}_0^{j+i}$  with at least one of the rows of  $\mathbf{G}_0^{j+i}, \mathbf{G}_1^{j+i}, \dots, \mathbf{G}_{K-1}^{j+i}$  and is therefore likely to attain each of its  $2^n$  values with equal probability as well and independently with the other code vectors (this fact is dealt in details in [10, Sect. 5.1]). (b) is true since we switched into looser conditions inside each event in the probability term. In (c) we used the union bound. In (d) we used the fact that observing  $\delta(\overline{\mathbf{V}}_j, \mathbf{Y}_j)$ , when summing up over all of  $\mathbf{V}_j$ 's possible values, is equivalent to observing  $\delta(\mathbf{V}_j, \mathbf{Y}_j)$  (since in both cases, each of the  $2^N$  values of the vector is covered by the summation). In (e) we used the Bhattacharyya bound for the pairwise error probability when using ML decision rule, and (f) is true since the channel is memoryless.

We proved that the probability that other code segment would be preferred by the minimax decoder over the correct segment, when averaged over the ensemble of time-varying

convolutional codes, is upper bounded by twice the bound achieved for ML decoder in [10]. Thus, it is exponentially of the same order. The subsequent steps in deriving an upper bound to the bit error exponent for rates  $R \leq R_{\theta,0}(Q^*)$  are identical to that of ML decoder (see [10, Sect. 5.1]) and the final result is the same.

Therefore, it was proved that when using the minimax decoder, the achievable exponent for bit error probability is no less than when the channel parameter is known and the ML decoder is used. The same error exponent was proved to be achievable for rates up to  $R_{\theta,0}(Q^*)$ .

In order to extend the average upper bound for the bit error probability to rates higher than  $R_{\theta,0}(Q^*)$ , we will use a slightly different technique.

First, we upper bound  $\pi_{l,\theta}(j)$ , the probability that a branch in the minimax based decoding path will occur by any one of the other possible paths, starting at node  $j$  and reemerging after  $K + l$  branches. We should observe, as mentioned in (45), that the number of such diverging paths satisfies  $M \leq (2^b - 1)2^{bl}$ . The code segments associated with these  $M$  incorrect paths will be denoted by  $\mathbf{V}_j^{(1)}, \dots, \mathbf{V}_j^{(M)}$ , respectively.

$$\pi_{l,\theta}(j) =$$

$$\begin{aligned} &= \Pr \left\{ \exists 1 \leq i \leq M : \rho(\mathbf{V}_j^{(i)}, \mathbf{Y}_j) \leq \rho(\mathbf{V}_j, \mathbf{Y}_j) \right\} \\ &= \Pr \left\{ \exists 1 \leq i \leq M : \min \left\{ \delta(\mathbf{V}_j^{(i)}, \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j^{(i)}, \mathbf{Y}_j) \right\} \leq \min \left\{ \delta(\mathbf{V}_j, \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \right\} \\ &= \Pr \left\{ \exists 1 \leq i \leq M : \delta(\mathbf{V}_j^{(i)}, \mathbf{Y}_j) \leq \min \left\{ \delta(\mathbf{V}_j, \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \bigcup \right. \\ &\quad \left. \delta(\overline{\mathbf{V}}_j^{(i)}, \mathbf{Y}_j) \leq \min \left\{ \delta(\mathbf{V}_j, \mathbf{Y}_j), 1 - \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \right\} \\ &\stackrel{(a)}{\leq} \Pr \left\{ \exists 1 \leq i \leq M : \delta(\mathbf{V}_j^{(i)}, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j) \bigcup \delta(\overline{\mathbf{V}}_j^{(i)}, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \\ &\stackrel{(b)}{\leq} \Pr \left\{ \exists 1 \leq i \leq M : \delta(\mathbf{V}_j^{(i)}, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} + \\ &\quad \Pr \left\{ \exists 1 \leq i \leq M : \delta(\overline{\mathbf{V}}_j^{(i)}, \mathbf{Y}_j) \leq \delta(\mathbf{V}_j, \mathbf{Y}_j) \right\} \\ &\stackrel{(c)}{\leq} \sum_{\mathbf{Y}_j} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \left[ \sum_{i=1}^M P_\theta(\mathbf{Y}_j | \mathbf{V}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho + \sum_{\mathbf{Y}_j} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \left[ \sum_{i=1}^M P_\theta(\mathbf{Y}_j | \overline{\mathbf{V}}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho, \quad (47) \end{aligned}$$

where (a) is true since we increased the right terms of the two inequalities, and thus increased the probability for union of these two events, in (b), the union bound was used, and in (c), we used the Gallager bound for the error probability when using the ML decision rule. This error was used for each of the two error probabilities.

We now move to upper bound the average of  $\pi_{l,\theta}(j)$  over the ensemble of time-varying convolutional codes:

$$\overline{\pi_{l,\theta}(j)} =$$

$$\begin{aligned}
&\stackrel{(a)}{=} \sum_{\mathbf{V}_j} 2^{-N} \sum_{\mathbf{V}_j^{(1)}} \dots \sum_{\mathbf{V}_j^{(M)}} 2^{-NM} \pi_{l,\theta}(j) \\
&\stackrel{(b)}{\leq} \sum_{\mathbf{Y}_j} \sum_{\mathbf{V}_j} 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \sum_{\mathbf{V}_j^{(1)}} \dots \sum_{\mathbf{V}_j^{(M)}} 2^{-NM} \\
&\quad \left[ \left[ \sum_{i=1}^M P_\theta(\mathbf{Y}_j | \mathbf{V}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho + \left[ \sum_{i=1}^M P_\theta(\mathbf{Y}_j | \overline{\mathbf{V}}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho \right] \\
&\stackrel{(c)}{=} 2 \cdot \sum_{\mathbf{Y}_j} \sum_{\mathbf{V}_j} 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \sum_{\mathbf{V}_j^{(1)}} \dots \sum_{\mathbf{V}_j^{(M)}} 2^{-NM} \left[ \sum_{i=1}^M P_\theta(\mathbf{Y}_j | \mathbf{V}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho \\
&\stackrel{(d)}{\leq} 2 \cdot \sum_{\mathbf{Y}_j} \sum_{\mathbf{V}_j} 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \left[ \sum_{i=1}^M \sum_{\mathbf{V}_j^{(1)}} \dots \sum_{\mathbf{V}_j^{(M)}} 2^{-NM} P_\theta(\mathbf{Y}_j | \mathbf{V}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho, \quad 0 \leq \rho \leq 1 \\
&\stackrel{(e)}{=} 2 \cdot \sum_{\mathbf{Y}_j} \sum_{\mathbf{V}_j} 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \left[ \sum_{i=1}^M \sum_{\mathbf{V}_j^{(i)}} 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho, \quad 0 \leq \rho \leq 1 \\
&\stackrel{(f)}{\leq} 2 \cdot (2^b - 1) 2^{bl\rho} \sum_{\mathbf{Y}_j} \sum_{\mathbf{V}_j} 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \left[ \sum_{\mathbf{V}_j^{(i)}} 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j^{(i)})^{\frac{1}{1+\rho}} \right]^\rho, \quad 0 \leq \rho \leq 1 \\
&= 2 \cdot (2^b - 1) 2^{bl\rho} \sum_{\mathbf{Y}_j} \sum_{\mathbf{V}_j} \left[ \sum 2^{-N} P_\theta(\mathbf{Y}_j | \mathbf{V}_j)^{\frac{1}{1+\rho}} \right]^{1+\rho}, \quad 0 \leq \rho \leq 1 \\
&\stackrel{(g)}{=} 2 \cdot (2^b - 1) 2^{bl\rho} \left\{ \sum_y \left[ \sum_v 2^{-N} P_\theta(y|v)^{\frac{1}{1+\rho}} \right]^{1+\rho} \right\}^N, \quad 0 \leq \rho \leq 1 \\
&\doteq (2^b - 1) 2^{bl\rho} e^{-(K+l)nE_{\theta,0}(\rho, \{\frac{1}{2}, \frac{1}{2}\})}, \quad 0 \leq \rho \leq 1, \tag{48}
\end{aligned}$$

where

$$E_{\theta,0}(\rho, Q^*) = -\ln \sum_y \left[ \sum_v 2^{-N} P_\theta(y|v)^{\frac{1}{1+\rho}} \right]^{1+\rho}.$$

In (a), we sum over all possible code vectors associated with the different paths in the trellis. As explained earlier, each code vector can attain all of its  $2^N$  values equiprobably and independently with the other code vectors. In (b), we used the result from (47). (c) is true since examining  $P_\theta(\mathbf{Y}_j | \overline{\mathbf{V}}_j^{(i)})$ ,  $1 \leq i \leq M$ , when summing up over all of  $\mathbf{V}_j^{(1)}, \dots, \mathbf{V}_j^{(M)}$  possible values, is equivalent to the examination of  $P_\theta(\mathbf{Y}_j | \mathbf{V}_j^{(i)})$ ,  $1 \leq i \leq M$ . In (d), we



bound ourselves to  $0 \leq \rho \leq 1$  and use Jensen's inequality. (e) is true since for a fixed  $i$ ,  $P_\theta(\mathbf{Y}_j|\mathbf{V}_j^{(i)})$  depends only on  $\mathbf{V}_j^{(i)}$ , and is enumerated for the  $2^{N(M-1)}$  possibilities of  $\mathbf{V}_j^{(1)}, \dots, \mathbf{V}_j^{(i-1)}, \mathbf{V}_j^{(i+1)}, \dots, \mathbf{V}_j^{(M)}$ . In (f), we upper bound  $M$  by  $(2^b - 1)2^{bl}$ , and (g) is true since the BSC is memoryless.

As in the above proof for rates up to  $R_{\theta,0}(Q^*)$ , the subsequent steps in deriving an upper bound to the bit error exponent for rates  $R_{\theta,0}(Q^*) \leq R \leq C_\theta$  for the minimax decoder are identical to that of ML decoder (see [10, Sect. 5.1]) and the final result is the same. This completes the proof that the achievable exponent for bit error probability of the minimax decoder is equal to that of the ML decoder, for all rates up to capacity.

## A. Appendix

### A.1 Proof of eq. (12) for ensembles of Linear and Systematic Linear Codes

In this section, we examine the performance of the minimax decoding rule with respect to uniform i.i.d. random coding over ensembles of linear codes and systematic linear codes. We will prove that for a family of BIOS channels, the same single-letter formula for the lower bound to the achievable fraction  $\xi^*$  is obtained, with uniform i.i.d. random coding distribution  $Q^* = \left\{\frac{1}{2}, \frac{1}{2}\right\}$  (i.e.  $\Delta^*(P) = \ln 2 - H(P)$ ).

Using Gallager's techniques, we first upper bound the decoding error probability given that the  $m$ -th message was sent for a given  $\theta$  in the following way:

$$\begin{aligned}
P_{E_m}(\Omega|\theta) &= \sum_{\mathbf{y} \in \mathcal{Y}^N} P_\theta(\mathbf{y}|\mathbf{v}_m) 1\left\{\exists m' \neq m : \max_{\theta' \in \Theta_N} f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y}) \geq \max_{\theta'' \in \Theta_N} f_{\theta''}(\mathbf{v}_m, \mathbf{y})\right\} \\
&= \sum_{\mathbf{y} \in \mathcal{Y}^N} P_\theta(\mathbf{y}|\mathbf{v}_m) 1\left\{\exists \theta', \exists m' \neq m : f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y}) \geq \max_{\theta'' \in \Theta_N} f_{\theta''}(\mathbf{v}_m, \mathbf{y})\right\} \\
&\stackrel{(a)}{=} \sum_{\mathbf{y} \in \mathcal{Y}^N} P_\theta(\mathbf{y}|\mathbf{v}_m) \max_{\theta' \in \Theta_N} 1\left\{\exists m' \neq m : f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y}) \geq \max_{\theta'' \in \Theta_N} f_{\theta''}(\mathbf{v}_m, \mathbf{y})\right\} \\
&= \sum_{\mathbf{y} \in \mathcal{Y}^N} P_\theta(\mathbf{y}|\mathbf{v}_m) \max_{\theta' \in \Theta_N} 1\left\{\exists m' \neq m : \frac{f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})}{\max_{\theta'' \in \Theta_N} f_{\theta''}(\mathbf{v}_m, \mathbf{y})} \geq 1\right\} \\
&\stackrel{(b)}{\leq} \sum_{\mathbf{y} \in \mathcal{Y}^N} P_\theta(\mathbf{y}|\mathbf{v}_m) \max_{\theta' \in \Theta_N} \min_{\substack{\lambda \geq 0 \\ \rho \geq 0}} \left[ \sum_{m' \neq m} \left( \frac{e^{N f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})}}{\max_{\theta'' \in \Theta_N} e^{N f_{\theta''}(\mathbf{v}_m, \mathbf{y})}} \right)^\lambda \right]^\rho, \quad (49)
\end{aligned}$$

where  $\rho \geq 0$  and  $\lambda \geq 0$  are free parameters.

(a) is true since if we denote with  $A(\theta)$  an event dependent on  $\theta \in \Theta_N$ , and denote with  $C$  a constant, then

$$1\{\exists \theta \in \Theta_N : A(\theta) > C\} = \max_{\theta \in \Theta_N} 1\{A(\theta) > C\}.$$

(b) is true since if we denote with  $f_1(m')$  and  $f_2(m)$  two non-negative functions of  $m'$  and  $m$  respectively, then (using Gallager's technique)

$$1 \left\{ \exists m' \neq m : \frac{f_1(m')}{f_2(m)} \geq 1 \right\} \leq \min_{\substack{\lambda \geq 0 \\ \rho \geq 0}} \left[ \sum_{m' \neq m} \left( \frac{e^{f_1(m')}}{e^{f_2(m)}} \right)^\lambda \right]^\rho.$$

Based on (49), we now develop an upper bound to the minimax criterion related to a specific linear code (i.e., specific values of  $\mathbf{G}$  and  $\mathbf{v}_0$ , thus denoted by  $S_N(\mathbf{v}_0, \mathbf{G})$ ):

$$\begin{aligned} S_N(\mathbf{v}_0, \mathbf{G}) &= \max_{\theta \in \Theta} \left\{ \frac{P_E(\Omega|\theta)}{[\bar{P}_E^*(\theta)]^\xi} \right\} = \max_{\theta \in \Theta} \left\{ \frac{P_E(\Omega|\theta)}{e^{-N\xi E_r^*(\theta)}} \right\} \\ &\leq \max_{\theta \in \Theta} \left\{ \frac{1}{M} \sum_{m=0}^{M-1} \sum_{\mathbf{y} \in \mathcal{Y}^N} e^{N\xi E_r^*(\theta)} P_\theta(\mathbf{y}|\mathbf{v}_m) \max_{\theta' \in \Theta_N} \min_{\substack{\lambda \geq 0 \\ \rho \geq 0}} \right. \\ &\quad \left. \left( \sum_{m' \neq m} \left[ \frac{e^{Nf_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})}}{\max_{\theta'' \in \Theta_N} e^{Nf_{\theta''}(\mathbf{v}_m, \mathbf{y})}} \right]^\lambda \right)^\rho \right\} \\ &= \max_{\theta \in \Theta} \left\{ \frac{1}{M} \sum_{m=0}^{M-1} \sum_{\mathbf{y} \in \mathcal{Y}^N} e^{N \cdot f_\theta(\mathbf{v}_m, \mathbf{y})} \max_{\theta' \in \Theta_N} \min_{\substack{\lambda \geq 0 \\ \rho \geq 0}} \right. \\ &\quad \left. \left( \sum_{m' \neq m} \left[ \frac{e^{Nf_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})}}{\max_{\theta'' \in \Theta_N} e^{Nf_{\theta''}(\mathbf{v}_m, \mathbf{y})}} \right]^\lambda \right)^\rho \right\} \\ &\stackrel{(a)}{\leq} \frac{1}{M} \sum_{m=0}^{M-1} \sum_{\mathbf{y} \in \mathcal{Y}^N} \left( \max_{\theta \in \Theta_N} e^{Nf_\theta(\mathbf{v}_m, \mathbf{y})} \right) \max_{\theta' \in \Theta_N} \min_{\substack{\lambda \geq 0 \\ \rho \geq 0}} \left\{ \frac{\left[ \sum_{m' \neq m} e^{N\lambda f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})} \right]^\rho}{\left( \max_{\theta'' \in \Theta_N} e^{Nf_{\theta''}(\mathbf{v}_m, \mathbf{y})} \right)^{\lambda\rho}} \right\} \\ &\stackrel{(b)}{\leq} \frac{1}{M} \sum_{m=0}^{M-1} \sum_{\mathbf{y} \in \mathcal{Y}^N} \max_{\theta' \in \Theta_N} \min_{\substack{\lambda \geq 0 \\ \rho \geq 0}} \left\{ \left[ \max_{\theta \in \Theta_N} e^{Nf_\theta(\mathbf{v}_m, \mathbf{y})} \right]^{1-\lambda\rho} \left[ \sum_{m' \neq m} e^{N\lambda f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})} \right]^\rho \right\} \\ &\stackrel{(c)}{\leq} \frac{1}{M} \sum_{m=0}^{M-1} \sum_{\mathbf{y} \in \mathcal{Y}^N} \max_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \max_{\theta \in \Theta_N} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_m, \mathbf{y})} \left[ \sum_{m' \neq m} e^{N\lambda f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})} \right]^\rho \right\} \\ &\stackrel{(d)}{=} \frac{1}{M} \sum_{m=0}^{M-1} \sum_{\mathbf{y} \in \mathcal{Y}^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_m, \mathbf{y})} \left[ \sum_{m' \neq m} e^{N\lambda f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})} \right]^\rho \right\}. \end{aligned} \tag{50}$$

The passages (a)–(d) are explained as follows: In (a) we used the fact that the maximum of an expectation is no greater than the expectation of the maximum and changed the maximization of  $\theta$  to be over  $\Theta_N$ . (b) is true since  $\theta$  and  $\theta''$  maximize two identical expressions, and therefore can be united. In (c) we restricted the range of the optimization to  $1 - \lambda\rho \geq 0$

( $\Rightarrow \lambda \leq 1/\rho$ ). In (d) we used the fact that for given  $\mathbf{v}_m, \mathbf{y}$  and  $\theta'$

$$\begin{aligned} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \max_{\theta \in \Theta_N} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_m, \mathbf{y})} \left[ \sum_{m' \neq m} e^{N\lambda f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})} \right]^\rho \right\} = \\ \max_{\theta \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_m, \mathbf{y})} \left[ \sum_{m' \neq m} e^{N\lambda f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})} \right]^\rho \right\}. \end{aligned} \quad (51)$$

This interchange between the minimization over  $\lambda$  and  $\rho$  and the maximization over  $\theta$  is justified in the Appendix, Section A.2.

Prior to deriving the single-letter formula for the lower bound to  $\xi^*$ , we first present the following claim:

**Lemma 2** *When a linear code is used for a BIOS channel and minimax decoding is used, the error probability for the  $m$ -th message is equal for all  $0 \leq m \leq M-1$ .*

This lemma is proved in Section A.4 of the Appendix.

Based on this observation, we can assume, without loss of generality, that  $\mathbf{u}_0 = \mathbf{0}$  was transmitted, and then the upper bound to  $S_N$  can be expressed as:

$$S_N(\mathbf{v}_0, \mathbf{G}) \leq \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ \sum_{m=1}^{M-1} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \quad (52)$$

In the following subsections, we will use the same technique to derive two upper bounds on the minimax criterion, one for the ensemble of linear codes and one for the ensemble of systematic linear codes.

## Linear Codes

By averaging  $S_N$  over the ensemble of linear codes:

$$\begin{aligned} \bar{S}_N &\stackrel{(a)}{=} 2^{-(K+1)N} \sum_{\mathbf{v}_0, \mathbf{G}} S_N(\mathbf{v}_0, \mathbf{G}) \\ &\leq 2^{-(K+1)N} \sum_{\mathbf{v}_0, \mathbf{G}} \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ \sum_{m=1}^{M-1} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \\ &\stackrel{(b)}{\leq} 2^{-(K+1)N} \sum_{\mathbf{v}_0, \mathbf{G}} \sum_{\mathbf{y} \in Y^N} \sum_{\theta \in \Theta_N} \sum_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ \sum_{m=1}^{M-1} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \end{aligned} \quad (53)$$

$$\begin{aligned}
& \stackrel{(c)}{\leq} \sum_{\mathbf{y} \in Y^N} \sum_{\theta \in \Theta_N} \sum_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ 2^{-(K+1)N} \sum_{\mathbf{v}_0, \mathbf{G}} \left( e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ \sum_{m=1}^{M-1} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right) \right\} \\
& \stackrel{(d)}{\leq} |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ 2^{-N} \sum_{\mathbf{v}_0} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} 2^{-KN} \sum_{\mathbf{G}} \left[ \sum_{m=1}^{M-1} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \tag{54} \\
& \stackrel{(e)}{\leq} |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ 2^{-N} \sum_{\mathbf{v}_0} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ 2^{-KN} \sum_{\mathbf{G}} \sum_{m=1}^{M-1} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \\
& = |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ 2^{-N} \sum_{\mathbf{v}_0} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ 2^{-KN} \sum_{m=1}^{M-1} \sum_{\mathbf{G}} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \\
& \stackrel{(f)}{=} |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ 2^{-N} \sum_{\mathbf{v}_0} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ (M-1) 2^{-N} \sum_{\mathbf{v}} e^{N\lambda f_{\theta'}(\mathbf{v}, \mathbf{y})} \right]^\rho \right\} \\
& \leq |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ M^\rho \left[ 2^{-N} \sum_{\mathbf{v}} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}, \mathbf{y})} \right] \cdot \left[ 2^{-N} \sum_{\mathbf{v}'} e^{N\lambda f_{\theta'}(\mathbf{v}', \mathbf{y})} \right]^\rho \right\}, \tag{55}
\end{aligned}$$

where the steps (a)–(f) are as follows: The equality in (a) is obtained by averaging over  $2^{(K+1)N}$  equiprobable values of  $\mathbf{v}_0$  and  $\mathbf{G}$ . (b) and (d) follow from the fact that for a non-negative function  $f(\theta)$ , non-negative function  $f(\theta)$ ,

$$\max_{\theta \in \Theta_N} f(\theta) \leq \sum_{\theta \in \Theta_N} f(\theta) \leq |\Theta_N| \cdot \max_{\theta \in \Theta_N} f(\theta). \tag{56}$$

(c) is true since an expectation of a minimum is upper-bounded by the minimum of the expectation. In (e), we limit the optimization over  $\rho$  to  $0 \leq \rho \leq 1$  and use Jensen's

inequality. In (f), we used the following equivalence for the two inner summations:

$$\sum_{m=1}^{M-1} \sum_{\mathbf{G}} e^{N\lambda f_{\theta'}(\mathbf{v}, \mathbf{m}, \mathbf{y})} = (M-1)2^{(K-1)N} \sum_{\mathbf{v}} e^{N\lambda f_{\theta'}(\mathbf{v}, \mathbf{y})}. \quad (57)$$

This equivalence is proved in Section A.5 of the Appendix.

From (19), we conclude that the term inside the summation in (55) is identical for all  $\mathbf{y}$ 's of the same type class. Thus, the summation can be conducted over types. Using (22), we continue to upper bound  $\bar{S}_N$  in the following way (note that the function  $A(\theta, \alpha, P_{\mathbf{x}\mathbf{y}})$  used here corresponds to a binary i.i.d. random coding distribution, as specified in (11)):

$$\begin{aligned} \bar{S}_N &\stackrel{(a)}{\leq} |\Theta_N|^2 \sum_{T_{\mathbf{y}}} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N\rho R} e^{N H_{\mathbf{y}}(Y)} e^{N \left[ (1-\lambda\rho)\xi E_r^*(\theta) - \min_{P_{\mathbf{x}|\mathbf{y}}} A(\theta, 1-\lambda\rho, P_{\mathbf{x}\mathbf{y}}) \right]} \right. \\ &\quad \left. e^{N \left[ \lambda\rho\xi E_r^*(\theta') - \rho \cdot \min_{P_{\mathbf{x}'|\mathbf{y}}} A(\theta', \lambda, P_{\mathbf{x}'\mathbf{y}}) \right]} \right\} \\ &= |\Theta_N|^2 \sum_{T_{\mathbf{y}}} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ \exp \left\{ N[\rho R + H_{\mathbf{y}}(Y) + (1-\lambda\rho)\xi E_r^*(\theta) - \min_{P_{\mathbf{x}|\mathbf{y}}} A(\theta, 1-\lambda\rho, P_{\mathbf{x}\mathbf{y}}) \right. \right. \\ &\quad \left. \left. + \lambda\rho\xi E_r^*(\theta') - \rho \cdot \min_{P_{\mathbf{x}'|\mathbf{y}}} A(\theta', \lambda, P_{\mathbf{x}'\mathbf{y}})] \right\} \right\} \\ &\stackrel{(b)}{\leq} |\Theta_N|^2 (N+1)^{|\mathcal{Y}|} \max_{P_{\mathbf{y}}} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ \exp \left\{ N[\rho R + H_{\mathbf{y}}(Y) + (1-\lambda\rho)\xi E_r^*(\theta) - \min_{P_{\mathbf{x}|\mathbf{y}}} A(\theta, 1-\lambda\rho, P_{\mathbf{x}\mathbf{y}}) \right. \right. \\ &\quad \left. \left. + \lambda\rho\xi E_r^*(\theta') - \rho \cdot \min_{P_{\mathbf{x}'|\mathbf{y}}} A(\theta', \lambda, P_{\mathbf{x}'\mathbf{y}})] \right\} \right\} \\ &= |\Theta_N|^2 (N+1)^{|\mathcal{Y}|} \cdot \exp \left\{ N \cdot \max_{P_{\mathbf{y}}} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \max_{P_{\mathbf{x}|\mathbf{y}}} \max_{P_{\mathbf{x}'|\mathbf{y}}} \right. \\ &\quad \left. [\rho R + H_{\mathbf{y}}(Y) + (1-\lambda\rho)\xi E_r^*(\theta) - A(\theta, 1-\lambda\rho, P_{\mathbf{x}\mathbf{y}}) \right. \\ &\quad \left. + \lambda\rho\xi E_r^*(\theta') - \rho \cdot A(\theta', \lambda, P_{\mathbf{x}'\mathbf{y}})] \right\}, \quad (58) \end{aligned}$$

where in (a) we upper bound  $|T_{\mathbf{y}}|$  by  $e^{N \cdot H_{\mathbf{y}}(Y)}$ , and in (b) we upper bound the summation of the functional over  $T_{\mathbf{y}}$  by the product of the maximal value (achieved by a specific distribution  $P_{\mathbf{y}}$ ) with  $(N+1)^{|\mathcal{Y}|}$ , which is an upper bound to the number of type classes  $\{T_{\mathbf{y}}\}$ .

As explained earlier, we seek the maximal  $\xi$  such that  $\overline{S}_N$  grows sub-exponentially with  $N$ . To this end, we can ignore the factor  $|\Theta_N|^2 (N+1)^{|\mathcal{Y}|}$  in (58), as it grows polynomially with  $N$ . Moreover, as mentioned in Section V, the optimizations can be conducted over continuous distributions and over the entire parameter space,  $\Theta$ . Thus, a maximal  $\xi$  is sought, such that (using (9)):

$$\max_{P_Y} \max_{\theta, \theta' \in \Theta} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \max_{P_{X|Y}} \max_{P_{X'|Y}} \left[ \rho R + (1-\lambda\rho)\xi E_r^*(\theta) + \lambda\rho\xi E_r^*(\theta') - B(\theta, \theta', P_Y, P_{X|Y}, P_{X'|Y}, \lambda, \rho) \right] \leq 0. \quad (59)$$

An equivalent condition to (59) is

$$\forall P_Y, \forall \theta, \theta' \in \Theta, \exists 0 \leq \rho \leq 1, 0 \leq \lambda \leq 1/\rho: \quad \forall P_{X|Y}, \forall P_{X'|Y}$$

$$\rho R + (1-\lambda\rho)\xi E_r^*(\theta) + \lambda\rho\xi E_r^*(\theta') - B(\theta, \theta', P_Y, P_{X|Y}, P_{X'|Y}, \lambda, \rho) \leq 0$$

or,

$$\forall P_Y, \forall \theta, \theta' \in \Theta, \exists 0 \leq \rho \leq 1, 0 \leq \lambda \leq 1/\rho: \quad \forall P_{X|Y}, \forall P_{X'|Y}$$

$$\xi \leq \frac{B(\theta, \theta', P_Y, P_{X|Y}, P_{X'|Y}, \lambda, \rho) - \rho R}{(1-\lambda\rho) \cdot E_r^*(\theta) + \lambda\rho \cdot E_r^*(\theta')}.$$

Consequently, for ensembles of linear codes and BIOS channels, the lower bound to  $\xi^*$  is the same as in (12), with a uniform i.i.d. random coding distribution,  $Q^* = \left\{ \frac{1}{2}, \frac{1}{2} \right\}$ .

### Systematic Linear Codes

A similar technique will be used now to achieve identical results for the ensemble of systematic linear codes.

By averaging  $S_N$  over this ensemble:

$$\begin{aligned} \overline{S}_N &\stackrel{(a)}{=} 2^{-K(N-K)} 2^{-N} \sum_{\tilde{\mathbf{G}}} \sum_{\mathbf{v}_0} S_N(\mathbf{v}_0, \mathbf{G}) \\ &\leq 2^{-K(N-K)} 2^{-N} \sum_{\tilde{\mathbf{G}}} \sum_{\mathbf{v}_0} \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \\ &\quad \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ \sum_{m=1}^{M-1} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \\ &\stackrel{(b)}{\leq} |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \\ &\quad \left\{ 2^{-N} \sum_{\mathbf{v}_0} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ 2^{-K(N-K)} \sum_{m=1}^{M-1} \sum_{\tilde{\mathbf{G}}} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} \right]^\rho \right\} \end{aligned}$$

$$\begin{aligned}
&\stackrel{(c)}{=} |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ 2^{-N} \sum_{\mathbf{v}_0} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ 2^{-(N-K)} \sum_{\mathbf{v}} e^{N\lambda f_{\theta'}(\mathbf{v}, \mathbf{y})} \right]^\rho \right\} \\
&\stackrel{(d)}{=} |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ 2^{-N} \sum_{\mathbf{v}_0} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_0, \mathbf{y})} \left[ M 2^{-N} \sum_{\mathbf{v}} e^{N\lambda f_{\theta'}(\mathbf{v}, \mathbf{y})} \right]^\rho \right\} \\
&= |\Theta_N|^2 \sum_{\mathbf{y} \in Y^N} \max_{\theta \in \Theta_N} \max_{\theta' \in \Theta_N} \min_{\substack{0 \leq \rho \leq 1 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ M^\rho \right. \\
&\quad \left. \left[ 2^{-N} \sum_{\mathbf{v}} e^{N(1-\lambda\rho)f_\theta(\mathbf{v}, \mathbf{y})} \right] \left[ 2^{-N} \sum_{\mathbf{v}'} e^{N\lambda f_{\theta'}(\mathbf{v}', \mathbf{y})} \right]^\rho \right\}
\end{aligned} \tag{60}$$

The equality in (a) is obtained by averaging over  $2^{K(N-K)}$  and  $2^N$  equiprobable values of  $\mathbf{v}_0$  and  $\tilde{\mathbf{G}}$  (the non-systematic part of  $\mathbf{G}$ ), respectively. (b) is obtained by taking identical steps as done for ensemble of linear codes in the previous subsection (see the inequalities between (53) and (54)). In (c), we used the following equivalence for the two inner summations:

$$\sum_{m=1}^{M-1} \sum_{\tilde{\mathbf{G}}} e^{N\lambda f_{\theta'}(\mathbf{v}_m, \mathbf{y})} = 2^{(K-1)(N-K)} \sum_{\mathbf{v}} e^{N\lambda f_{\theta'}(\mathbf{v}, \mathbf{y})}. \tag{61}$$

This equivalence is proved in Section A.6 of the Appendix. In (d), we used the equality  $M = 2^K$ .

Finally, the upper bound to  $\bar{\mathcal{S}}_N$  achieved in (60) is identical to the one related to ensembles of linear codes (see (55)), and therefore the final lower bound to  $\xi^*$  for the case of systematic linear codes is also identical to (12) with uniform i.i.d. random coding distribution,  $Q^* = \left\{ \frac{1}{2}, \frac{1}{2} \right\}$ .

## A.2 Proof of eq. (26) and eq. (51)

Let  $\theta^*$  maximize  $f_\theta(\mathbf{v}_m, \mathbf{y})$ , and let  $F(\lambda, \rho)$  be a nonnegative function. Then,

$$\begin{aligned}
&\min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \max_{\theta \in \Theta_N} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_m, \mathbf{y})} \cdot F(\lambda, \rho) \right\} = \\
&= \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N(1-\lambda\rho)f_{\theta^*}(\mathbf{v}_m, \mathbf{y})} \cdot F(\lambda, \rho) \right\} \\
&\stackrel{(a)}{\leq} \max_{\theta \in \Theta_N} \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_m, \mathbf{y})} \cdot F(\lambda, \rho) \right\} \\
&\leq \min_{\substack{\rho \geq 0 \\ 0 \leq \lambda \leq 1/\rho}} \max_{\theta \in \Theta_N} \left\{ e^{N(1-\lambda\rho)f_\theta(\mathbf{v}_m, \mathbf{y})} \cdot F(\lambda, \rho) \right\},
\end{aligned} \tag{62}$$

where (a) is true since the value of the function for a specific  $\theta^*$  in  $\Theta_N$  is always upper-bounded by the maximization of the function over  $\theta \in \Theta_N$ . Thus, all inequalities must be achieved with equalities.

### A.3 Proof of eq. (22)

For  $\alpha \in \Re$  and  $\mathbf{y} \in \mathcal{Y}^N$ , we exponentially evaluate  $E[e^{N\alpha f_\theta(\mathbf{x}, \mathbf{y})}]$ , where the average is calculated over the ensemble of random coding distribution of the form:  $Q_N(\mathbf{x}) = \frac{Q_N(T\mathbf{x})}{|T\mathbf{x}|}$ , (as described in (22)):

$$\begin{aligned}
E[e^{N\alpha f_\theta(\mathbf{X}, \mathbf{y})}] &= \sum_{\mathbf{x} \in \mathcal{X}^N} Q_N(\mathbf{x}) e^{N\alpha f_\theta(\mathbf{x}, \mathbf{y})} \\
&\stackrel{(a)}{=} \sum_{T\mathbf{x}|\mathbf{y} \subset \mathcal{X}^N} |T\mathbf{x}|\mathbf{y}| Q_N(\mathbf{x}) e^{N\alpha f_\theta(\mathbf{x}, \mathbf{y})} \\
&= \sum_{T\mathbf{x}|\mathbf{y} \subset \mathcal{X}^N} |T\mathbf{x}|\mathbf{y}| \frac{e^{-N\Delta_N^*(P\mathbf{x})}}{|T\mathbf{x}|} e^{N\alpha f_\theta(\mathbf{x}, \mathbf{y})} \\
&\stackrel{(b)}{=} \sum_{T\mathbf{x}|\mathbf{y} \subset \mathcal{X}^N} |T\mathbf{x}|\mathbf{y}| \frac{e^{-N(\Delta^*(P\mathbf{x}) - \tilde{\epsilon}_N)}}{|T\mathbf{x}|} e^{N\alpha f_\theta(\mathbf{x}, \mathbf{y})}, \tag{63}
\end{aligned}$$

where  $\tilde{\epsilon}_N \rightarrow 0$  as  $N \rightarrow \infty$  independently of  $P\mathbf{x}$ .

We should note that (a) is true since  $f_\theta(\mathbf{x}, \mathbf{y})$  depends on  $\mathbf{x}$  and  $\mathbf{y}$  only via their joint empirical distribution and the summation can be conducted over types instead, and since the average is calculated for a given  $\mathbf{y}$ , we sum over  $T\mathbf{x}|\mathbf{y}$ . In (b) we used the convergence assumption for the random coding distributions withing the class  $\mathcal{Q}$ .

Thus, we continue to evaluate  $E[e^{N\alpha f_\theta(\mathbf{x}, \mathbf{y})}]$  as follows:

$$\begin{aligned}
E[e^{N\alpha f_\theta(\mathbf{x}, \mathbf{y})}] &\stackrel{(a)}{=} \sum_{T\mathbf{x}|\mathbf{y} \subset \mathcal{X}^N} \exp\left\{N[-I\mathbf{x}\mathbf{y}(X; Y) - \Delta^*(P\mathbf{x}) + \alpha f_\theta(\mathbf{x}, \mathbf{y})]\right\} \\
&\stackrel{(b)}{=} \exp\left\{N \cdot \max_{P\mathbf{x}|\mathbf{y}} [-I\mathbf{x}\mathbf{y}(X; Y) - \Delta^*(P\mathbf{x}) + \alpha f_\theta(\mathbf{x}, \mathbf{y})]\right\} \\
&\stackrel{(c)}{=} \exp\left\{N \cdot \max_{P\mathbf{x}|\mathbf{y}} [-I\mathbf{x}\mathbf{y}(X; Y) - \Delta^*(P\mathbf{x}) \right. \\
&\quad \left. + \alpha \hat{E}\mathbf{x}\mathbf{y} \ln P_\theta(Y|X) + \alpha \xi E_r^*(\theta)]\right\} \\
&= \exp\left\{N[\alpha \xi E_r^*(\theta) - \min_{P\mathbf{x}|\mathbf{y}} \{I\mathbf{x}\mathbf{y}(X; Y) \right. \\
&\quad \left. + \Delta^*(P\mathbf{x}) - \alpha \hat{E}\mathbf{x}\mathbf{y} \ln P_\theta(Y|X)\}]\right\} \\
&= e^{N[\alpha \xi E_r^*(\theta) - \min_{P\mathbf{x}|\mathbf{y}} A(\theta, \alpha, P\mathbf{x}\mathbf{y})]}, \tag{64}
\end{aligned}$$



where in (a), we used the facts that  $|T_{\mathbf{x}|\mathbf{y}}| \doteq e^{N \cdot H \mathbf{x} \mathbf{y}^{(X|Y)}}$  and  $|T_{\mathbf{x}}| \doteq e^{N \cdot H \mathbf{x}^{(X)}}$ . (b) is true since the summation of the functional over  $T_{\mathbf{x}|\mathbf{y}} \subset \mathcal{X}^N$  is lower bounded by its maximal value (achieved by a specific distribution  $P_{\mathbf{x}|\mathbf{y}}$ ), and upper bounded by the product of its maximal value with  $(N+1)^{|\mathcal{X}||\mathcal{Y}|}$ . In (c), we expressed the minimax metric in terms of the joint empirical distribution as described in (19).

#### A.4 Proof of Lemma 2

In this section, we prove that when a linear code is used for a BIOS channel and the minimax decision rule is used (denoted by  $\Omega$ ), the error probability for the  $m$ -th message (of length  $N$ ),  $\mathbf{v}_m = (v_{m0}, \dots, v_{m(N-1)})$ , is the same for all  $m$ , that is,

$$P_{E_m}(\Omega|\theta) = P_E(\Omega|\theta) \quad \text{for } 0 \leq m \leq M-1. \quad (65)$$

Considering a binary input channel, we denote the channel crossover probabilities for a single letter as  $P_\theta(y|v=0) \triangleq P_{\theta,0}(y)$  and  $P_\theta(y|v=1) \triangleq P_{\theta,1}(y)$ .

If the channel is also output symmetric then,

$$P_{\theta,1}(y) = P_{\theta,0}(-y), \quad \forall y \in \mathcal{Y}$$

The error probability for the  $m$ -th message using minimax decoding is:

$$\begin{aligned} P_{E_m}(\Omega|\theta) &= \sum_{\mathbf{y} \in \Lambda_m^c} P_\theta(\mathbf{y}|\mathbf{v}_m) \\ &= \sum_{\mathbf{y} \in \Lambda_m^c} \prod_{n:v_{mn}=0} P_{\theta,0}(y_n) \prod_{n:v_{mn}=1} P_{\theta,1}(y_n) \\ &= \sum_{\mathbf{y} \in \Lambda_m^c} \prod_{n:v_{mn}=0} P_{\theta,0}(y_n) \prod_{n:v_{mn}=1} P_{\theta,0}(-y_n), \end{aligned} \quad (66)$$

where

$$\begin{aligned} \Lambda_m^c &= \left\{ \mathbf{y} : \max_{\theta'} \left\{ \frac{1}{N} \ln P_{\theta'}(\mathbf{y}|\mathbf{v}_{m'}) + \xi E_r^*(\theta') \right\} \geq \max_{\theta''} \left\{ \frac{1}{N} \ln P_{\theta''}(\mathbf{y}|\mathbf{v}_m) + \xi E_r^*(\theta'') \right\}, \right. \\ &\quad \left. \text{for some } m' \neq m \right\} \\ &= \left\{ \mathbf{y} : \max_{\theta'} \left\{ \sum_{n=0}^{N-1} \ln P_{\theta'}(y_n|v_{m'n}) + N \xi E_r^*(\theta') \right\} \geq \right. \\ &\quad \left. \max_{\theta''} \left\{ \sum_{n=0}^{N-1} \ln P_{\theta''}(y_n|v_{mn}) + N \xi E_r^*(\theta'') \right\}, \quad \text{for some } m' \neq m \right\} \end{aligned}$$

$$\begin{aligned}
&= \left\{ \mathbf{y} : \max_{\theta'} \left\{ \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=0}} \ln P_{\theta',0}(y_t) + \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=1}} \ln P_{\theta',1}(y_t) + \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=0}} \ln P_{\theta',0}(y_t) + \right. \\
&\quad \left. \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=1}} \ln P_{\theta',1}(y_t) + N\xi E_r^*(\theta') \right\} \geq \\
&\quad \max_{\theta''} \left\{ \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=0}} \ln P_{\theta'',0}(y_t) + \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=1}} \ln P_{\theta'',0}(y_t) + \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=0}} \ln P_{\theta'',1}(y_t) + \right. \\
&\quad \left. \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=1}} \ln P_{\theta'',1}(y_t) + N\xi E_r^*(\theta'') \right\}, \quad \text{for some } m' \neq m \Big\} \\
&= \left\{ \mathbf{y} : \max_{\theta'} \left\{ \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=0}} \ln P_{\theta',0}(y_t) + \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=1}} \ln P_{\theta',0}(-y_t) + \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=0}} \ln P_{\theta',0}(y_t) + \right. \\
&\quad \left. \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=1}} \ln P_{\theta',0}(-y_t) + N\xi E_r^*(\theta') \right\} \geq \\
&\quad \max_{\theta''} \left\{ \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=0}} \ln P_{\theta'',0}(y_t) + \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=1}} \ln P_{\theta'',0}(y_t) + \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=0}} \ln P_{\theta'',0}(-y_t) + \right. \\
&\quad \left. \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=1}} \ln P_{\theta'',0}(-y_t) + N\xi E_r^*(\theta'') \right\}, \quad \text{for some } m' \neq m \Big\}. \quad (67)
\end{aligned}$$

Using the following transformation to dummy variables

$$z_n = \begin{cases} y_n, & \forall n : v_{mn}=0 \\ -y_n, & \forall n : v_{mn}=1 \end{cases}$$

we get that

$$\begin{aligned}
P_{E_m}(f|\theta) &= \sum_{\mathbf{z} \in \Lambda_m^c} \prod_{n: v_{mn}=0} P_{\theta,0}(z_n) \prod_{n: v_{mn}=1} P_{\theta,0}(z_n) \\
&= \sum_{\mathbf{z} \in \Lambda_m^c} \prod_{n=0}^{N-1} P_{\theta,0}(z_n), \quad (68)
\end{aligned}$$

where

$$\begin{aligned}
\Lambda_m^c &= \left\{ \mathbf{z} : \max_{\theta'} \left\{ \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=0}} \ln P_{\theta',0}(z_t) + \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=1}} \ln P_{\theta',0}(-z_t) + \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=0}} \ln P_{\theta',0}(-z_t) + \right. \\
&\quad \left. \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=1}} \ln P_{\theta',0}(z_t) + N\xi E_r^*(\theta') \right\} \geq \\
&\quad \max_{\theta''} \left\{ \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=0}} \ln P_{\theta'',0}(z_t) + \sum_{\substack{t: v_{mt}=0 \\ v_{m't}=1}} \ln P_{\theta'',0}(z_t) + \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=0}} \ln P_{\theta'',0}(z_t) + \right. \\
&\quad \left. \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=1}} \ln P_{\theta'',0}(z_t) + N\xi E_r^*(\theta'') \right\} \Big\}
\end{aligned}$$

$$\begin{aligned}
& \left. \sum_{\substack{t: v_{mt}=1 \\ v_{m't}=1}} \ln P_{\theta'',0}(z_t) + N\xi E_r^*(\theta'') \right\}, \quad \text{for some } m' \neq m \Bigg\} \\
&= \left\{ \mathbf{z} : \max_{\theta'} \left\{ \sum_{p: v_{mp}=v_{m'p}} \ln P_{\theta',0}(z_p) + \sum_{q: v_{mq} \neq v_{m'q}} \ln P_{\theta',0}(-z_q) + N\xi E_r^*(\theta') \right\} \geq \right. \\
& \quad \max_{\theta''} \left\{ \sum_{p: v_{mp}=v_{m'p}} \ln P_{\theta'',0}(z_p) + \sum_{q: v_{mq} \neq v_{m'q}} \ln P_{\theta'',0}(z_q) + N\xi E_r^*(\theta'') \right\}, \\
& \quad \left. \text{for some } m' \neq m \right\}. \tag{69}
\end{aligned}$$

Now, on the one hand, (68) and (69) describe  $P_{E_m}(f|\theta)$  and  $\Lambda_m^c$ , respectively, for each  $0 \leq m \leq M-1$ . On the other hand, we should note that the terms for  $P_{E_0}(f|\theta)$  and  $\Lambda_0^c$  (describing the case where  $\mathbf{v}_0 = \mathbf{0}$  is transmitted) are obtained by assigning  $m = 0$  in (66) and (67). By doing that, the result terms coincide with (68) and (69), respectively (which, as mentioned before, correspond to the  $m$ -th message). This observation completes the proof.

### A.5 Proof of eq. (57)

First, by the way of constructing the linear code, we know that:

$$\mathbf{v}_{m'} = \mathbf{u}_{m'} \mathbf{G} \oplus \mathbf{v}_0, \quad \forall 0 \leq m' \leq M-1 \tag{70}$$

Since  $1 \leq m' \leq M-1$  implies  $\mathbf{u}_{m'} \neq \mathbf{0}$ , then for each information vector in this set there is at least one index  $i$  for which  $u_{m'i} = 1$ . Consequently, the construction of each code vector  $\mathbf{v}_{m'}$ ,  $1 \leq m' \leq M-1$ , can be written in the following way:

$$\mathbf{v}_{m'} = \mathbf{u}_{m'} \mathbf{G} \oplus \mathbf{v}_0 = \mathbf{g}_i \oplus \left[ \sum_{j \neq i} u_{m'j} \mathbf{g}_j \right] \oplus \mathbf{v}_0,$$

where  $\mathbf{g}_i$  stands for the  $i$ -th row in  $\mathbf{G}$ .

Therefore:

$$\begin{aligned}
\sum_{m'=1}^{M-1} \sum_{\mathbf{G}} e^{N\lambda f_{\theta'}(\mathbf{v}_{m'}, \mathbf{y})} &= \sum_{m'=1}^{M-1} \sum_{\mathbf{G} \setminus \mathbf{g}_i} \sum_{\mathbf{g}_i} \exp \left\{ N f_{\theta'}(\mathbf{g}_i \oplus \left[ \sum_{j \neq i} u_{m'j} \mathbf{g}_j \right] \oplus \mathbf{v}_0, \mathbf{y}) \lambda \right\} \\
&\stackrel{(a)}{=} \sum_{m'=1}^{M-1} \sum_{\mathbf{G} \setminus \mathbf{g}_i} \sum_{\mathbf{v}} e^{N f_{\theta'}(\mathbf{v}, \mathbf{y}) \lambda} \\
&= \sum_{m'=1}^{M-1} 2^{(K-1)N} \sum_{\mathbf{v}} e^{N f_{\theta'}(\mathbf{v}, \mathbf{y}) \lambda}
\end{aligned}$$

$$= (M-1)2^{(K-1)N} \sum_{\mathbf{v}} e^{Nf_{\theta'}(\mathbf{v}; \mathbf{y})\lambda}, \quad (71)$$

where (a) is true since for fixed values of  $m'$ ,  $\mathbf{G} \setminus \mathbf{g}_i$  (in the outer summations) and  $\mathbf{v}_0$ , the row vector, which is denoted by  $\mathbf{v}$ , is fixed, causing  $\mathbf{g}_i$  to sum up over all the binary vectors of length  $N$ .

## A.6 Proof of eq. (61)

In this section, we prove the equality, which is given in (61), and used in (60).

First, by the way of constructing a systematic linear code:

$$\begin{aligned} \mathbf{v}_{m'} &= \mathbf{u}_{m'} \mathbf{G} \oplus \mathbf{v}_0 \\ &= \left[ \mathbf{u}_{m'}; \sum_{i=1}^K u_{m'i} \tilde{\mathbf{g}}_i \right] \oplus \mathbf{v}_0 \\ &= \left[ \mathbf{u}_{m'}; \overbrace{0 \dots 0}^{N-K} \right] \oplus \left[ \overbrace{0 \dots 0}^K; \sum_{i=1}^K u_{m'i} \tilde{\mathbf{g}}_i \right] \oplus \mathbf{v}_0, \quad \forall 0 \leq m' \leq M-1, \end{aligned} \quad (72)$$

where  $\tilde{\mathbf{g}}_i$  stands for the  $i$ 'th row in  $\tilde{\mathbf{G}}$  (the non-systematic part of  $\mathbf{G}$ ).

We observe that for  $1 \leq m' \leq M-1$ ,  $\mathbf{u}_{m'} \neq \mathbf{0}$ . Thus, for each information vector in this set there's at least one index  $i$  for which  $u_{m'i} = 1$ . Consequently, the construction of each code vector  $\mathbf{v}_{m'}$ ,  $1 \leq m' \leq M-1$ , can be written in the following way:

$$\mathbf{v}_{m'} = \left[ \mathbf{u}_{m'}; \tilde{\mathbf{g}}_i \right] \oplus \left[ \overbrace{0 \dots 0}^K; \sum_{j \neq i} u_{m'j} \tilde{\mathbf{g}}_j \right] \oplus \mathbf{v}_0.$$

Therefore:

$$\begin{aligned} \sum_{m'=1}^{M-1} \sum_{\tilde{\mathbf{G}}} e^{Nf_{\theta'}(\mathbf{v}_{m'}; \mathbf{y})\lambda} &= \sum_{m'=1}^{M-1} \sum_{\tilde{\mathbf{G}} \setminus \tilde{\mathbf{g}}_i} \sum_{\tilde{\mathbf{g}}_i} \exp \left\{ Nf_{\theta'} \left( \left[ \mathbf{u}_{m'}; \tilde{\mathbf{g}}_i \right] \oplus \left[ 0 \dots 0; \sum_{j \neq i} u_{m'j} \tilde{\mathbf{g}}_j \right] \oplus \mathbf{v}_0, \mathbf{y} \right) \lambda \right\} \\ &\stackrel{(a)}{=} \sum_{\tilde{\mathbf{G}} \setminus \tilde{\mathbf{g}}_i} \sum_{m'=1}^{M-1} \sum_{\tilde{\mathbf{g}}_i} e^{Nf_{\theta'} \left( \left[ \mathbf{u}_{m'}; \tilde{\mathbf{g}}_i \right] \oplus \mathbf{v}, \mathbf{y} \right) \lambda} \\ &\stackrel{(b)}{\leq} \sum_{\tilde{\mathbf{G}} \setminus \tilde{\mathbf{g}}_i} \sum_{m'=0}^{M-1} \sum_{\tilde{\mathbf{g}}_i} e^{Nf_{\theta'} \left( \left[ \mathbf{u}_{m'}; \tilde{\mathbf{g}}_i \right] \oplus \mathbf{v}, \mathbf{y} \right) \lambda} \\ &\stackrel{(c)}{=} \sum_{\tilde{\mathbf{G}} \setminus \tilde{\mathbf{g}}_i} \sum_{\mathbf{v}} e^{Nf_{\theta'}(\mathbf{v}; \mathbf{y})\lambda} \\ &= 2^{(K-1)(N-K)} \sum_{\mathbf{v}} e^{Nf_{\theta'}(\mathbf{v}; \mathbf{y})\lambda}, \end{aligned} \quad (73)$$

where (a) is true since for fixed values of  $m'$ ,  $\tilde{\mathbf{G}} \setminus \tilde{\mathbf{g}}_i$  (in the outer summations) and  $\mathbf{v}_0$ , the row vector, which is denoted by  $\mathbf{v}$ , is fixed. In (b),  $m' = 0$  was added to the summation, and since the inner term in the summation is always non-negative the result cannot get smaller. (c) is true since for a fixed  $\mathbf{v}$ , summing up over  $0 \leq m' \leq M - 1$  and  $\tilde{\mathbf{g}}_i$  is equivalent to the summation over all the possibilities for a vector of length  $N$ .

## A.7 Equivalence between decision rules - $\Omega$ and $\Lambda$

In this section, we prove the equivalence between the minimax decision rule,  $\Omega$ , maximizing the metric  $f(\mathbf{x}, \mathbf{y})$  (as defined in (3)), and a decision rule  $\Lambda$ , minimizing  $\rho(\mathbf{x}, \mathbf{y})$  (as defined in (17)). We will prove that for a given output  $\mathbf{y} \in \mathcal{Y}$ , each  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  satisfy:

$$f(\mathbf{x}_1, \mathbf{y}) \geq f(\mathbf{x}_2, \mathbf{y}) \iff \rho(\mathbf{x}_1, \mathbf{y}) \leq \rho(\mathbf{x}_2, \mathbf{y}). \quad (74)$$

First, we should note that  $f(\mathbf{x}, \mathbf{y})$  satisfies:

$$\begin{aligned} f(\mathbf{x}, \mathbf{y}) &= \max_{0 \leq \theta \leq 1} f_\theta(\mathbf{x}, \mathbf{y}) \\ &= \max_{0 \leq \theta \leq 1} \left\{ \frac{1}{N} [\ln P_\theta(\mathbf{y}|\mathbf{x}) + N\xi E_r^*(\theta)] \right\} \\ &\stackrel{(a)}{=} \max_{0 \leq \theta \leq 1} \left\{ \frac{1}{N} [d(\mathbf{x}, \mathbf{y}) \ln \theta + (N - d(\mathbf{x}, \mathbf{y})) \ln (1 - \theta) + N\xi E_r^*(\theta)] \right\} \\ &= \max_{0 \leq \theta \leq 1} \left\{ \delta(\mathbf{x}, \mathbf{y}) \ln \theta + (1 - \delta(\mathbf{x}, \mathbf{y})) \ln (1 - \theta) + \xi E_r^*(\theta) \right\} \\ &= \max_{0 \leq \theta \leq 1} f_\theta(\delta(\mathbf{x}, \mathbf{y})), \quad 0 \leq \delta(\mathbf{x}, \mathbf{y}) \leq 1 \\ &= f(\delta(\mathbf{x}, \mathbf{y})), \quad 0 \leq \delta(\mathbf{x}, \mathbf{y}) \leq 1. \end{aligned} \quad (75)$$

In (a), we used the following representation for the BSC transition probability:

$$P_\theta(\mathbf{y}|\mathbf{x}) = \theta^{d(\mathbf{x}, \mathbf{y})} (1 - \theta)^{N - d(\mathbf{x}, \mathbf{y})}.$$

We conclude that the value of  $f(\mathbf{x}, \mathbf{y})$  is equal for all code vectors with the same (normalized) Hamming distance from  $\mathbf{y}$ , and therefore can be defined as  $f(\delta(\mathbf{x}, \mathbf{y}))$ ,  $0 \leq \delta(\mathbf{x}, \mathbf{y}) \leq 1$ .

Next, we now prove that  $f(\mathbf{x}, \mathbf{y})$  has the same value for a code vector  $\mathbf{x}$  and its complement,  $\bar{\mathbf{x}}$ :

$$\begin{aligned} f(\bar{\mathbf{x}}, \mathbf{y}) &= f(\delta(\bar{\mathbf{x}}, \mathbf{y})) \\ &= f(1 - \delta(\mathbf{x}, \mathbf{y})) \\ &= \max_{0 \leq \theta \leq 1} f_\theta(1 - \delta(\mathbf{x}, \mathbf{y})) \end{aligned}$$

$$\begin{aligned}
&= \max_{0 \leq \theta \leq 1} \left\{ (1 - \delta(\mathbf{x}, \mathbf{y})) \ln \theta + \delta(\mathbf{x}, \mathbf{y}) \ln (1 - \theta) + \xi E_r^*(\theta) \right\} \\
&\stackrel{(a)}{=} \max_{0 \leq \tilde{\theta} \leq 1} \left\{ (1 - \delta(\mathbf{x}, \mathbf{y})) \ln (1 - \tilde{\theta}) + \delta(\mathbf{x}, \mathbf{y}) \ln \tilde{\theta} + \xi E_r^*(1 - \tilde{\theta}) \right\} \\
&\stackrel{(b)}{=} \max_{0 \leq \tilde{\theta} \leq 1} \left\{ (1 - \delta(\mathbf{x}, \mathbf{y})) \ln (1 - \tilde{\theta}) + \delta(\mathbf{x}, \mathbf{y}) \ln \tilde{\theta} + \xi E_r^*(\tilde{\theta}) \right\} \\
&= \max_{0 \leq \tilde{\theta} \leq 1} f_{\tilde{\theta}}(\delta(\mathbf{x}, \mathbf{y})) \\
&= f(\delta(\mathbf{x}, \mathbf{y})) \\
&= f(\mathbf{x}, \mathbf{y}).
\end{aligned} \tag{76}$$

In (a), we changed the variable in the maximization,  $\tilde{\theta} = 1 - \theta$ , and (b) is true since for the BSC model the ML error exponent,  $E_r^*(\theta)$ , is symmetric around  $\theta = \frac{1}{2}$  (see (42)).

Using the fact that both  $f(\delta(\mathbf{x}, \mathbf{y}))$  and  $\rho(\delta(\mathbf{x}, \mathbf{y}))$  are equal for  $\delta(\mathbf{x}, \mathbf{y})$  and  $1 - \delta(\mathbf{x}, \mathbf{y})$ , it is sufficient to prove (74) for  $\mathbf{x}_1$  and  $\mathbf{x}_2$  satisfying  $\delta(\mathbf{x}_1, \mathbf{y}) \leq \frac{1}{2}$  and  $\delta(\mathbf{x}_2, \mathbf{y}) \leq \frac{1}{2}$  (and thus  $\rho(\mathbf{x}_1, \mathbf{y}) = \delta(\mathbf{x}_1, \mathbf{y})$ ,  $\rho(\mathbf{x}_2, \mathbf{y}) = \delta(\mathbf{x}_2, \mathbf{y})$ ).

In the rest of the proof, we will denote  $\delta(\mathbf{x}_1, \mathbf{y}) \triangleq \delta_1$ ,  $\delta(\mathbf{x}_2, \mathbf{y}) \triangleq \delta_2$ . It is therefore sufficient to show that

$$f(\delta_1) \geq f(\delta_2) \iff 0 \leq \delta_1 \leq \delta_2 \leq \frac{1}{2}. \tag{77}$$

This equivalence will be shown in two steps:

First, we note that  $0 \leq \delta_1 \leq \delta_2 \leq \frac{1}{2}$  satisfy that  $\forall 0 \leq \theta \leq \frac{1}{2}$ :

$$\delta_1 \ln \left( \frac{\theta}{1 - \theta} \right) \geq \delta_2 \ln \left( \frac{\theta}{1 - \theta} \right). \tag{78}$$

By adding  $\ln(1 - \theta) + \xi E_r^*(\theta)$  to both sides of (78) we get:

$$\delta_1 \ln \left( \frac{\theta}{1 - \theta} \right) + \ln(1 - \theta) + \xi E_r^*(\theta) \geq \delta_2 \ln \left( \frac{\theta}{1 - \theta} \right) + \ln(1 - \theta) + \xi E_r^*(\theta) \tag{79}$$

or

$$\delta_1 \ln \theta + (1 - \delta_1) \ln(1 - \theta) + \xi E_r^*(\theta) \geq \delta_2 \ln \theta + (1 - \delta_2) \ln(1 - \theta) + \xi E_r^*(\theta). \tag{80}$$

This inequality is true for the values of  $\theta$ , which maximize the both sides of (80). i.e.:

$$\begin{aligned}
&\max_{0 \leq \theta \leq \frac{1}{2}} \{ \delta_1 \ln \theta + (1 - \delta_1) \ln(1 - \theta) + \xi E_r^*(\theta) \} \geq \\
&\max_{0 \leq \theta \leq \frac{1}{2}} \{ \delta_2 \ln \theta + (1 - \delta_2) \ln(1 - \theta) + \xi E_r^*(\theta) \}
\end{aligned} \tag{81}$$

or

$$\max_{0 \leq \theta \leq \frac{1}{2}} f_{\theta}(\delta_1) \geq \max_{0 \leq \theta \leq \frac{1}{2}} f_{\theta}(\delta_2). \tag{82}$$

In order to complete the proof, one must broaden the maximization ranges over  $\theta$  in (82) into  $0 \leq \theta \leq 1$ . In order to justify that this broadening is possible, we present the following observation:

Each  $0 \leq \delta \leq \frac{1}{2}$  satisfy that  $\forall \frac{1}{2} \leq \theta \leq 1$ :

$$\delta \ln \left( \frac{\theta}{1-\theta} \right) \leq (1-\delta) \ln \left( \frac{\theta}{1-\theta} \right). \quad (83)$$

By adding  $\ln(1-\theta) + \xi E_r^*(\theta)$  to both sides of (83) we get:

$$\delta \ln \left( \frac{\theta}{1-\theta} \right) + \ln(1-\theta) + \xi E_r^*(\theta) \leq (1-\delta) \ln \left( \frac{\theta}{1-\theta} \right) + \ln(1-\theta) + \xi E_r^*(\theta) \quad (84)$$

or

$$\delta \ln \theta + (1-\delta) \ln(1-\theta) + \xi E_r^*(\theta) \leq \delta \ln(1-\theta) + (1-\delta) \ln \theta + \xi E_r^*(\theta). \quad (85)$$

Using the fact that for the BSC model the ML error exponent,  $E_r^*(\theta)$ , is symmetric around  $\theta = \frac{1}{2}$  (see (42)), we can rewrite (85) as:

$$\delta \ln \theta + (1-\delta) \ln(1-\theta) + \xi E_r^*(\theta) \leq \delta \ln(1-\theta) + (1-\delta) \ln \theta + \xi E_r^*(1-\theta) \quad (86)$$

or

$$f_\theta(\delta(\mathbf{x}, \mathbf{y})) \leq f_{1-\theta}(\delta(\mathbf{x}, \mathbf{y})). \quad (87)$$

The meaning of (87) is that when  $0 \leq \delta \leq \frac{1}{2}$ , for each  $\frac{1}{2} \leq \theta \leq 1$ ,  $f_\theta(\delta)$  is always upper bounded by  $f_{1-\theta}(\delta)$  where  $0 \leq 1-\theta \leq \frac{1}{2}$ . Thus, maximization of  $f_\theta(\delta)$  over  $0 \leq \theta \leq 1$  is obviously accomplished by  $\theta$  in  $\left[0, \frac{1}{2}\right]$ .

Therefore, (82) finally becomes:

$$\max_{0 \leq \theta \leq 1} f_\theta(\delta_1) \geq \max_{0 \leq \theta \leq 1} f_\theta(\delta_2) \quad (88)$$

thus,

$$0 \leq \delta_1 \leq \delta_2 \leq \frac{1}{2} \Leftrightarrow f(\delta_1) \geq f(\delta_2), \quad (89)$$

and the proof is complete.

## References

- [1] I. Csiszár, “Linear Codes for Sources and Source Networks: Error Exponents, Universal Coding,” *IEEE Trans. Inform. Theory*, vol. IT-28, pp. 585-592, July 1982.
- [2] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, Academic Press 1981.
- [3] M. Feder and A. Lapidoth, “Universal Decoders for Channels with Memory,” *IEEE Trans. Inform. Theory*, vol. IT-44, no. 5, pp. 1726–1745, September 1998.
- [4] M. Feder and N. Merhav, “Universal Composite Hypothesis Testing - A Competitive Minimax Approach,” *IEEE Trans. Inform. Theory*, vol. IT-48, pp. 1504-1517, June 2002.
- [5] V. D. Goppa, “Nonprobabilistic Mutual Information Without Memory,” *Probl. Cont. Information Theory*, Vol. 4, pp. 97-102, 1975.
- [6] A. Lapidoth and J. Ziv, “On the Universality of the LZ-based Decoding Algorithm,” *IEEE Trans. Inform. Theory*, vol. IT-44, no. 5, pp. 1746–1755, September 1998.
- [7] N. Merhav and M. Feder, “Minimax Universal Decoding with an Erasure Option,” *IEEE Trans. Inform. Theory*, vol. IT-53, no. 5, pp. 1664–1675, May 2007.
- [8] N. Shulman, “Communication over an Unknown Channel via Common Broadcasting,” Ph.D. dissertation, Tel Aviv University, July 2003.
- [9] A. Somekh-Baruch and N. Merhav, “Achievable error exponents for the private fingerprinting game,” *IEEE Trans. Inform. Theory*, vol. IT-53, no. 5, pp. 1827–1838, May 2007.
- [10] A. J. Viterbi and J. K. Omura, *Principles of Digital Communication and Coding*. McGraw-Hill, 1979.
- [11] O. Zeitouni, J. Ziv, and N. Merhav, “When is the Generalized Likelihood Ratio Test Optimal?” *IEEE Trans. Inform. Theory*, vol. IT-38, no. 5, pp. 1597–1602, September 1992.



- [12] J. Ziv, “Universal Decoding for Finite-State Channels,” *IEEE Trans. Inform. Theory*, vol. IT-31, no. 4, pp. 453–460, July 1985.